

Online Visual Place Recognition via Saliency Re-identification

Han Wang, Chen Wang, and Lihua Xie

Abstract—As an essential component of visual simultaneous localization and mapping (SLAM), place recognition is crucial for robot navigation and autonomous driving. Existing methods often formulate visual place recognition as feature matching, which is computationally expensive for many robotic applications with limited computing power, *e.g.*, autonomous driving and cleaning robot. Inspired by the fact that human beings always recognize a place by remembering salient regions or landmarks that are more attractive or interesting than others, we formulate visual place recognition as saliency re-identification. In the meanwhile, we propose to perform both saliency detection and re-identification in frequency domain, in which all operations become element-wise. The experiments show that our proposed method achieves competitive accuracy and much higher speed than the state-of-the-art feature-based methods. The proposed method is open-sourced and available at <https://github.com/wh200720041/SRLCD.git>.

I. INTRODUCTION

Visual place recognition, also known as loop closure detection, is the task to identify repetitive places or landmarks [1], [2] during autonomous navigation. It is of great importance for generating drift-free maps in simultaneous localization and mapping (SLAM). In robot navigation, trajectory estimation often comes with drifts due to the sensor imperfection and environmental variation [3]. Without the capability of loop closure, estimated robot pose inevitably deviates from its true position. This subsequently leads to an unreliable localization and also increases the computational expense due to repetitive registration of landmarks [4]. Therefore, to enable a drift-free localization, a well-structured SLAM system usually requires visual place recognition to associate current pose with historical frames.

Computational cost is one of the key bottlenecks in visual place recognition [5], [6]. For example, in a SLAM system, the query times of incoming data grows as more places are visited and registered into the database (or map), which can subsequently slow down the system or cause memory overflow in the long run. Most of existing works down-scale an image (or a place) into lower dimensional features to improve computational efficiency. A popular strategy is to re-formulate the visual place recognition problem as the retrieval of the local hand-crafted features such as ORB [7] and BRISK [8]. It has been widely used in many applications due to its high accuracy, *e.g.*, brute-force (BF) matching on

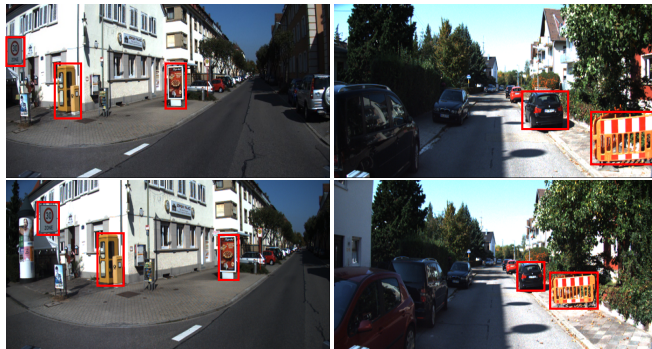


Fig. 1: Instead of matching feature points, we formulate place recognition as saliency re-identification. Both saliency detection and retrieval are performed in frequency domain to take advantage of efficiency of element-wise operation. Different from general object detection, saliency detection extracts areas that are visually appealing. This figure shows that the salient regions are re-identified when the places are re-visited (samples from KITTI dataset [12]).

local descriptors has achieved satisfactory results on public dataset [9]. However, to guarantee a high accuracy, feature-based approaches require the extraction and matching of thousands of descriptors for each image, which is computationally expensive, *e.g.*, iBoW-LCD [10] and FABMAP [11]. Moreover, discretization of local descriptors requires complicated offline training that is troublesome in practice.

In many real-time robotic systems such as unmanned aerial vehicles (UAVs) and micro-robots, computational capability is very limited due to power supply and payload constraints. Hence an effective place representation strategy is necessary.

Inspired by the fact that human beings always recognize a place by remembering objects or landmarks that are more appealing or interesting than others, not simply by remembering point features [13], we formulate place recognition as saliency re-identification, which is more natural and straightforward than feature matching. However, the definition of saliency is subjective and can be varied according to one’s psychological process [14] and environmental changes [15]. For simplicity, we follow the observation that human eyes are more sensitive to salient regions in an image with specific spectrum properties [16]. The salient regions can be either objects or visual distinct areas that stand out and enclosed by bounding boxes. In this paper, we find that both saliency detection and retrieval can be performed in frequency domain, which can further reduce the computational cost.

The work is supported by Delta-NTU Corporate Laboratory for Cyber-Physical Systems under the National Research Foundation Corporate Lab @ University Scheme.

Han Wang and Lihua Xie are with the School of Electrical and Electronic Engineering, Nanyang Technological University, 50 Nanyang Avenue, Singapore 639798. e-mail: {wang.han, elhxie}@ntu.edu.sg

Chen Wang is with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213-3890 USA. e-mail: chenwang@dr.com

In summary, the main contributions of this paper are:

- We propose a novel framework for visual place recognition for SLAM via saliency re-identification, which is coded in C++ and open sourced.
- We propose to detect and identify the salient regions in frequency domain by taking advantage of computationally efficient element-wise operation.
- The experiments show that our method is much faster than existing feature-based methods without reducing recall rate and precision. It does not require offline training or loading trained vocabulary.

This paper is organized as follows: Section II reviews the related works on visual place recognition. Section III describes the details of the proposed framework. Section IV shows experimental results and comparison with existing methods, followed by conclusion in Section V.

II. RELATED WORK

In this section, we review the most widely used feature-based methods, including both offline training-based approach and online training-based approach. We will also review the most recent deep learning-based methods.

A. Feature-based Methods

Although proposed early, feature matching is still the most popular strategy for visual place recognition, especially in robotics. One of the most classic approaches is the fast appearance-based mapping (FABMAP) [17]. In FABMAP, the visual vocabulary of SURF feature [18] is trained by hierarchical k-means clustering to discretize its high dimensional representation. Revisited places are then identified by matching feature distribution over the offline trained visual vocabulary. Similar idea is employed in DBoW2 [19], in which the binary descriptor ORB [7] is used instead. It achieves faster speed due to the binary feature representation. However, discretization of local descriptors requires complicated offline training for visual vocabulary that is troublesome in practice. For example, in FABMAP, it takes a few hours to train such visual vocabulary on a desktop computer. Moreover, the training data cannot contain any loop and have to be collected in the similar environment in order to achieve a good result.

Instead of building an offline trained visual vocabulary to discretize local feature space, some recent works introduce more efficient feature retrieval to make online training possible. For example, Emilio *et al.* [10] introduce a vocabulary maintenance strategy called dynamic island that groups similar images. The dynamic island identifies repetitive registration of same descriptor across multiple frames. The database is kept at small scale by removing those redundant descriptors. Schlegel *et al.* [9] introduce an online trained Hamming distance embedded binary search tree (HBST) to image retrieval which is much faster than traditional FLANN matching methods. However, the memory cost is huge since raw descriptors are used to build the incremental visual vocabulary tree. Applying dimensional reduction on local features is also an alternative solution.

Carrasco *et al.* [20] adopt the idea from image encoding to extract local descriptors. Hash coding is applied to the local features array collected from single image so that the extracted features are more compact. Gehrig *et al.* [21] directly apply principle component analysis (PCA) for BRISK [8] feature. Instead of searching on pre-trained visual vocabulary, k-nearest neighbor (K-NN) search is performed on the projected descriptors that achieves faster speed at the level of millisecond for each query. However, those methods require local feature extraction, which is computationally expensive.

B. Deep Learning-based Methods

The recent success of convolutional neural network (CNN) in computer vision [22] has triggered another research trend for visual place recognition. For example, in [23], a multi-scale feature encoding method is introduced by training two CNN architectures. The generated CNN features are viewpoint invariant, hence a large performance improvement is achieved. Inspired by the traditional image retrieval method, *i.e.*, vector of locally aggregated descriptors (VLAD), Relja *et al.* propose the NetVLAD [24] to learn CNN parameters in an end-to-end manner for place recognition. Hausler *et al.* [25] combine both feature-based methods and the CNN techniques. The query image is trained by combining sum of absolute difference (SAD), histogram of oriented gradients (HOG) [26], CNN spatial Max pooling and CNN spatial Arg-Max pooling. However, those CNN-based methods achieve high accuracy at the expense of a huge computational burden, which is sometimes infeasible for real-time systems.

III. METHODOLOGY

The proposed framework mainly consists of two components: saliency detection and saliency retrieval, both of which will be processed in frequency domain to take advantage of the efficiency of element-wise operation.

A. Saliency Detection

Existing methods on saliency analysis mainly focus on spectral analysis [13]. In particular, we introduce log-spectral residual approach to identify visual saliency which has been widely applied in practice [16]. Log-spectral analysis can be computed in an efficient element-wise manner, which can subsequently reduce the computational cost.

1) *Definition:* We present the basic procedure of log-spectral residual for saliency detection. Denote the 2-D FFT $\mathcal{F} : \mathbb{C}^{m \times n} \mapsto \mathbb{C}^{m \times n}$ as $\hat{\cdot}$, given an input image \mathbf{I} , its FFT can be expressed as $\hat{\mathbf{I}}$. We denote the magnitude spectrum and phase spectrum as $\mathcal{A}(\hat{\mathbf{I}})$ and $\mathcal{P}(\hat{\mathbf{I}})$, respectively, thus the log-spectral residual $\mathcal{R}(\hat{\mathbf{I}})$ can be defined as:

$$\mathcal{L}(\hat{\mathbf{I}}) = \log\left(\mathcal{A}(\hat{\mathbf{I}})\right), \quad (1a)$$

$$\mathcal{R}(\hat{\mathbf{I}}) = \exp\left(\mathcal{L}(\hat{\mathbf{I}}) - \mathbf{f}_{ave} * \mathcal{L}(\hat{\mathbf{I}})\right), \quad (1b)$$

where \mathbf{f}_{ave} is a normalized average filter of size $k \times k$ and $*$ denotes the cross-correlation. A saliency map \mathcal{M} can be

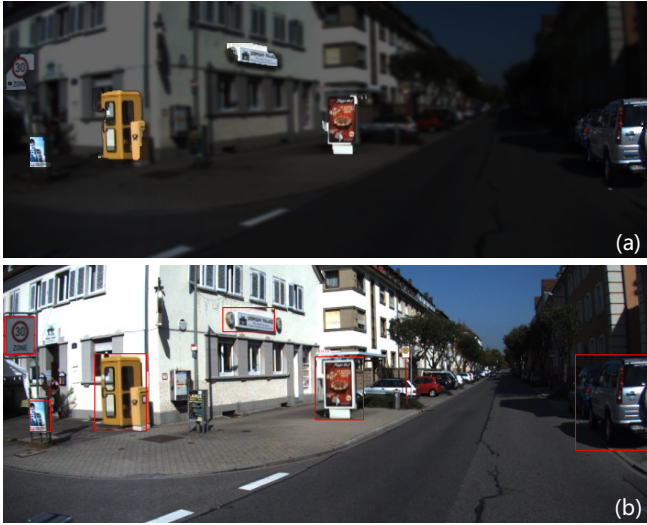


Fig. 2: An example of saliency detection from KITTI dataset. (a) shows the mask of the detected saliency map. (b) is the detected salient regions enclosed by red bounding boxes.

simply extracted by converting the residual back to spatial domain via inverse fast Fourier transform (IFFT):

$$\mathcal{M} = \mathcal{G}_\sigma * \mathcal{F}^{-1} \left(\mathcal{R}(\hat{\mathbf{I}}) \cdot \exp(\mathbf{j} \cdot \mathcal{P}(\hat{\mathbf{I}})) \right), \quad (2)$$

where \mathcal{G}_σ is a Gaussian filter and \mathbf{j} is the imaginary unit. This saliency map (2) implies the saliency distribution of the original image [27]. With such distribution, we can derive saliency regions $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbf{I}$ by simply taking pixel connectivity analysis on the saliency map \mathcal{M} .

The log-spectral residual approach provides a fast solution for saliency detection. However, not all extracted saliency information is suitable for place recognition, *e.g.*, the regions that are too small or too dark or bright. Therefore, we need extra rules to filter out the unqualified regions in order to improve the robustness.

It is observed that the detected saliency with high contrast and more edges is often unique and highly distinguishable, which is suitable for place recognition. Considering the computational cost, we adopt a discriminative strategy using contrast density ϕ_x and edge complexity ρ_x defined in (3).

$$\phi_x = \overline{(\mathbf{x} - \bar{\mathbf{x}})^2}, \quad (3a)$$

$$\rho_x = \frac{\mathcal{C}(\mathbf{x})}{\mathcal{S}(\mathbf{x})}, \quad (3b)$$

where $\mathbf{x} \in \mathbf{I}$ is the extracted salient region, $\bar{\mathbf{x}}$ is the mean pixel intensity, \mathcal{C} is Canny edge filter, and \mathcal{S} computes the total size of salient region (*i.e.*, pixel number). In experiment, we find that this strategy can remove most of less informative regions and even perform better than more complicated rules.

The output of log-spectrum residual algorithm is a mask with roughly object shapes. For better matching, we enclose those masks by minimum bounding boxes. An example of saliency detection is shown in Fig. 2. It can be seen that six distinct regions, including traffic signs, ad boards, phone

booths, and cars are extracted. Although it is inevitable that some dynamic objects, *e.g.*, cars, are extracted, they will be further removed by consistency verification in Section III-C.

2) *Property*: One of the reasons that we take the log-spectral residual for saliency detection is that the detected results are translation, rotation, and scale equivariant, which is crucial for visual place recognition. Recall that the corresponding FFT of translation-rotated image $\mathbf{I}_2(x, y) = \mathbf{I}_1(x \cos \theta + y \sin \theta - x_0, -x \sin \theta + y \cos \theta - y_0)$ is related by $\hat{\mathbf{I}}_2(u, v) = e^{-j2\pi(ux_0 + vy_0)} \cdot \hat{\mathbf{I}}_1(u \cos \theta + v \sin \theta, -u \sin \theta + v \cos \theta)$ [28], in which the factor $e^{-j2\pi(ux_0 + vy_0)}$ does not change the magnitude spectrum. This means that log-spectrum residual $\mathcal{R}(\hat{\mathbf{I}}_2)$ is just a rotation replica of $\mathcal{R}(\hat{\mathbf{I}}_1)$, resulting in a rotated saliency map \mathcal{M} and unchanged salient regions. This also holds for scale transforms. Suppose \mathbf{I}_2 is a scale transform of \mathbf{I}_1 with scale factor s , thus $\hat{\mathbf{I}}_2(u, v) = \frac{1}{s^2} \hat{\mathbf{I}}_1(u/s, v/s)$ and $\mathcal{L}(\hat{\mathbf{I}}_2(u, v)) = \mathcal{L}(\hat{\mathbf{I}}_1(u/s, v/s)) - 2 \log s$, meaning that the log-spectrum is scaled and the same salient region can be extracted. Those properties are crucial for robust saliency retrieval.

B. Saliency Retrieval

Although the saliency extraction is robust and consistent, the extracted salient regions can slightly differ from size, viewing angles, illumination, *etc.* Therefore, a transform-invariant encoding algorithm is necessary for saliency retrieval. Considering the computational cost, we introduce the kernel cross-correlator (KCC) [29], which is able to match two image areas directly in frequency domain and is invariant to affine transforms such as translation, rotation, and scale. More specifically, we compare the current salient region \mathbf{z} with the salient region \mathbf{x} from database. The purpose of training is to find an optimal correlator that is unique for each salient region. Then the correlator is used to examine the similarity of two salient regions in the retrieval stage.

1) *Definition*: Recall that the cross-correlation of two 2-D signals $\mathbf{g} = \mathbf{x} * \mathbf{h}$ becomes $\hat{\mathbf{g}} = \hat{\mathbf{x}} \odot \hat{\mathbf{h}}^*$ in frequency domain, where the operator $*$, \odot and \mathbf{h}^* represent cross-correlation, element-wise multiplication, and complex conjugate, respectively. Given correlation output \mathbf{g} , the kernel cross-correlation is defined as:

$$\hat{\mathbf{g}} = \hat{\kappa}_z(\mathbf{x}) \odot \hat{\mathbf{h}}^*, \quad (4)$$

where $\kappa_z(\mathbf{x})$ is a kernel matrix given by:

$$\kappa_z(\mathbf{x}) = \begin{bmatrix} \kappa(\mathbf{x}, \mathbf{z}_{11}) & \dots & \kappa(\mathbf{x}, \mathbf{z}_{1n}) \\ \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}, \mathbf{z}_{m1}) & \dots & \kappa(\mathbf{x}, \mathbf{z}_{mn}) \end{bmatrix}, \quad (5)$$

with $\kappa(\cdot, \cdot)$ being a kernel function. \mathbf{x} and \mathbf{z} are the candidate salient region and target salient region. \mathbf{z}_{ij} is the translational shift of \mathbf{z} with i pixel shifts in horizontal and j pixel shifts in vertical direction. For a salient region \mathbf{z} of size $m \times n$, there are $m \times n$ translational shifts in total.

Specifically, we use the Gaussian kernel $\kappa(\mathbf{x}, \mathbf{z}) = \phi(\|\mathbf{x} - \mathbf{z}\|^2)$. Let function $\Phi(\cdot)$ be the matrix form of Gaussian

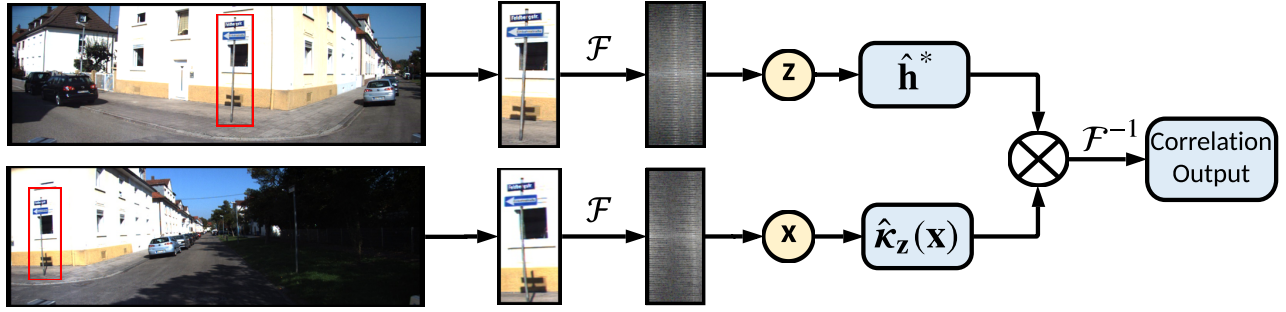


Fig. 3: The procedure of saliency retrieval. The same traffic sign is identified when a car re-visits the same place.

kernel function, where each element is a single Gaussian kernel $\phi(\cdot)$, then the kernel matrix $\kappa_{\mathbf{z}}(\mathbf{x})$ becomes:

$$\begin{aligned} \kappa_{\mathbf{z}}(\mathbf{x}) &= \begin{bmatrix} \phi(\|\mathbf{x} - \mathbf{z}_{11}\|^2) & \dots & \phi(\|\mathbf{x} - \mathbf{z}_{1n}\|^2) \\ \vdots & \ddots & \vdots \\ \phi(\|\mathbf{x} - \mathbf{z}_{m1}\|^2) & \dots & \phi(\|\mathbf{x} - \mathbf{z}_{mn}\|^2) \end{bmatrix} \quad (6) \\ &= \Phi(\|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 - 2 \cdot [\text{Tr}(\mathbf{x}^T \mathbf{z}_{ij})]_{mn}), \end{aligned}$$

where $\text{Tr}(\cdot)$ is the matrix trace and $[\text{Tr}(\mathbf{x}^T \mathbf{z}_{ij})]_{mn}$ denotes:

$$[\text{Tr}(\mathbf{x}^T \mathbf{z}_{ij})]_{mn} = \begin{bmatrix} \text{Tr}(\mathbf{x}^T \mathbf{z}_{11}) & \dots & \text{Tr}(\mathbf{x}^T \mathbf{z}_{1n}) \\ \vdots & \ddots & \vdots \\ \text{Tr}(\mathbf{x}^T \mathbf{z}_{m1}) & \dots & \text{Tr}(\mathbf{x}^T \mathbf{z}_{mn}) \end{bmatrix}, \quad (7)$$

From the 2-D correlation theory, we know that $\mathbf{x} * \mathbf{z} = [\text{Tr}(\mathbf{x}^T \mathbf{z}_{ij})]_{mn}$, so that we have

$$\kappa_{\mathbf{z}}(\mathbf{x}) = \Phi(\|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 - 2 \cdot \mathbf{x} * \mathbf{z}) \quad (8a)$$

$$= \Phi\left(\frac{\|\hat{\mathbf{x}}\|^2}{m \cdot n} + \frac{\|\hat{\mathbf{z}}\|^2}{m \cdot n} - 2 \cdot \mathcal{F}^{-1}(\hat{\mathbf{x}} \odot \hat{\mathbf{z}}^*)\right). \quad (8b)$$

This reduces the complexity $\mathcal{O}(N^2)$ of (6) to $\mathcal{O}(N \log N)$ of (8b), where $N = m \times n$ is the number of pixels in the salient region. This implies that KCC can also be calculated efficiently due to the element-wise operation in frequency domain [30].

2) *Training*: Given desired correlation output \mathbf{g} , we expect an optimal \mathbf{h} that maps the kernel matrix $\kappa_{\mathbf{z}}(\mathbf{x})$ to \mathbf{g} such that the sum of squared error (SSE) is minimized in frequency domain:

$$\min_{\hat{\mathbf{h}}^*} \|\hat{\kappa}_{\mathbf{z}}(\mathbf{x}) \odot \hat{\mathbf{h}}^* - \hat{\mathbf{g}}\|^2 + \lambda \|\hat{\mathbf{h}}^*\|^2. \quad (9)$$

It has a closed-form solution:

$$\hat{\mathbf{h}}^* = \frac{\hat{\mathbf{g}}}{\hat{\kappa}_{\mathbf{z}}(\mathbf{x}) + \lambda}, \quad (10)$$

where operator \div is the element-wise division. In the training stage, the optimal correlator \mathbf{h} is solved with $\mathbf{x} = \mathbf{z}$. Note that due to the efficiency of (10), we can obtain $\hat{\mathbf{h}}^*$ online, hence no pre-trained database is required.

3) *Retrieval*: As shown in [29], kernel cross-correlation is equivariant to affine transforms, *i.e.*, if input \mathbf{x} is transformed (translation, scale, and rotation), the output \mathbf{g} will be translated accordingly. Therefore, the transform of a test sample \mathbf{x} can be estimated by examining the translation of correlation output in (4). For simplicity, we set \mathbf{g} in (10) as

$$\mathbf{g}(x, y) = \begin{cases} 1, & \text{if } x = 0, y = 0 \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

Because of noise and distortion, we cannot always obtain a pure translational response of \mathbf{g} in (11). Therefore, we can take the maximum value in the test output as the similarity ζ between training and test samples, which will be invariant to the affine transforms of \mathbf{x} .

$$\zeta = \max \mathcal{F}^{-1}(\hat{\kappa}_{\mathbf{z}}(\mathbf{x}) \odot \hat{\mathbf{h}}^*). \quad (12)$$

The procedure of saliency retrieval using KCC is presented in Fig. 3, in which a loop closure candidate is detected.

C. Consistency Verification

As mentioned before, re-identification through saliency can improve the computational efficiency. However, the matched salient region may also cause false positive because the approach only identifies the similarity between query and candidate salient regions. Therefore, it is necessary to check consistency before closing the loop since any false positive will lead to localization and mapping failure easily.

The consistency verification is performed only on detected loop pairs by brute force matching using ORB features. Despite computationally expensive, matching on raw feature descriptor is able to remove most of the false positives produced by saliency re-identification, implying that traditional feature matching is a good complementary technique for saliency matching. We extract 1000 descriptors from current image and candidate image. A loop closure is determined when the number of matched descriptor pairs exceeds certain threshold. Fortunately, the loop closure does not occur frequently and the computational cost is less significant in comparison to the entire process of place recognition. In the practical experiment, it costs less than 5% of total computational expense so that such computational cost is worthy to guarantee the geometrical consistency.

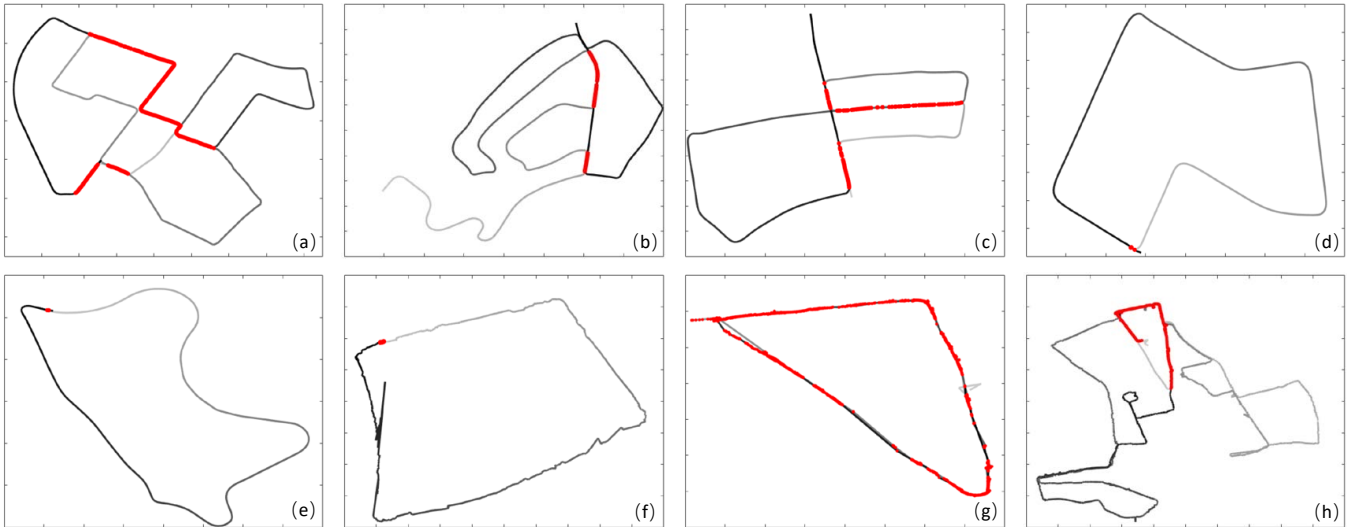


Fig. 4: Examples of place recognition on public dataset. The trajectory from GPS is plotted in from light grey to dark grey with time going. The loop closure reported by the proposed method is marked with red dot. (a, b, c, d, e) are collected from KITTI sequence 00, 02, 05, 07 and 09 respectively. (f, g, h) are collected from Oxford RobotCar sequence 2014_05_14_13_50_20, 2014_08_11_10_22_21 and 2014_12_02_15_30_08 respectively. Our method achieves 80% recall rate with no false positive.

IV. EXPERIMENTS

A. Setup and Metric

In our experiment, the proposed method and the compared methods are coded in C++ and tested on an *Intel i7-8700 3.2 GHz* CPU. The source codes of compared methods are also included in provided link. The parameter list used in this paper is presented in Table I.

Three metrics, *i.e.*, recall precision, recall rate, and computational cost, that are widely used to evaluate place recognition, are adopted in this paper. Specifically, recall precision is the percentage of success match; recall rate is the reported matches against total matches obtained from ground truth; computational cost refers to the processing time. To note that, for practical applications, recall precision is required to be as high as possible since false positive will result in false mapping easily. Computational cost is of great concern as most of robotic systems do not possess enough computational power. Recall rate is considered as satisfied if most of the loop places are identified.

B. Performance evaluation

The evaluation is performed based on various public datasets that are popularly used for loop closure detection, including KITTI [12], TUM [31], New College [32], City

Center [32], and Oxford RobotCar [33]. They include various scenarios consisting of indoor localization and autonomous driving, *etc.* The datasets contain large scale collections of up to 35,000 frames (RobotCar dataset), medium scale collections of 5,000 frames (KITTI dataset) and small scale collections of 1,000 frames (New College dataset). The details are listed in Table II. The ground truth of loop closure detection is collected based on GPS/Vicon information.

a) Results: Fig. 4 shows the examples of loop closure detection by the proposed approach. We plot the results from large datasets for illustration. The datasets include more than 60,000 images in total and the largest dataset contains around 35,000 images. As can be seen, most of re-visited places are identified and there is no false positive.

b) Precision-Recall Curve: The precision-recall curve on all datasets is shown in Fig. 5a. In particular, RobotCar 2014_05, 2014_06, 2014_08, and 2014_12 refers to RobotCar sequence 2014_05_14_13_50_20, 2014_06_26_09_53_12, 2014_08_11_10_22_21, and 2014_12_02_15_30_08, respectively. It is notable that our method achieves 100% precision over all datasets. In some challenging outdoor dynamic environments, *e.g.*, Oxford RobotCar dataset with 35,000 images, our approach achieves high recall rate up to 90%. More importantly, our approach is able to run up to 50Hz, which is

Parameter	Description	Value
σ	Standard deviation of Gaussian filter \mathcal{G}_σ	2.0
n	Size of average filter \mathbf{f}_{ave}	7
ϕ_{th}	Threshold of contrast density	58
ρ_{th}	Threshold of edge complexity	0.5
η_{th}	Threshold of KCC similarity test η	0.4

TABLE I: Parameter List.

Dataset	Description	Image Resolution
KITTI [12]	Outdoor, dynamic	370×1226
Freiburg3 [31]	Indoor, static	1280×1024
New College [32]	Outdoor, dynamic	640×480
City Center [32]	Outdoor, dynamic	640×480
RobotCar [33]	Outdoor, dynamic	1280×960

TABLE II: Details of Datasets.

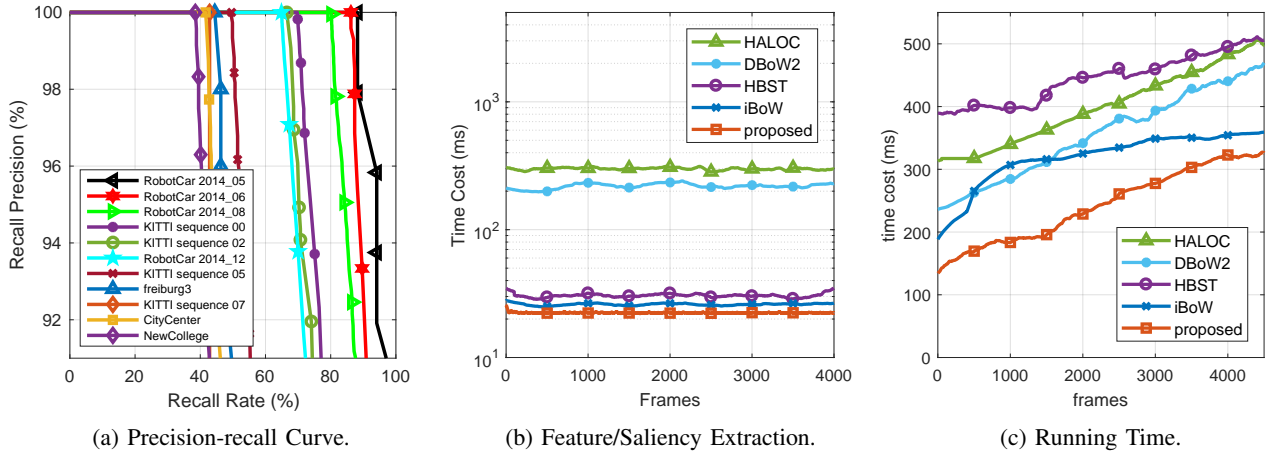


Fig. 5: (a) Precision-recall curve of the proposed method on different dataset. (b) Comparison of saliency extraction with feature extraction methods in KITTI sequence 00. (c) Time cost comparison of different approaches on KITTI sequence 00; Our approach achieves much faster speed than existing methods.

computationally cost-effective for many robotic applications.

C. Comparison

To further illustrate the efficiency of the proposed method, we compare the state-of-the-art feature-based methods such as iBoW [10], HALOC [20], DBoW2 [19], and HBST [9]. In particular, DBoW2 appears in the early stage and yet is still one of the most popular methods and has been used for ORB SLAM [34] and LDSO [35]. The rest methods are introduced in the recent years and they have been shown to be effective in loop closure detection. For all compared methods, one thousand local features, which is a typical recommended feature extraction number, are extracted for each frame. To note that, HBST is evaluated at quarter of original video size due to its huge memory consumption and time cost.

a) *Efficiency Evaluation*: For computational efficiency, we use the KITTI sequence 00 for example since the speed comparison is similar across different datasets. Fig. 5b shows

the performance of the saliency extraction compared to feature extraction from other methods. As mentioned before, saliency extraction is only of linear computational complexity in frequency domain, it is much faster than local feature extraction. In iBoW and HBST, raw ORB feature is used and hence the computational cost is also low. However, HALOC and DBoW2 require further dimension reduction hence the feature extraction time is higher. Fig. 5c shows the overall computational cost of respective methods. Our approach achieves the fastest speed among all the approaches, *i.e.*, a few times faster than DBoW2 and HALOC.

b) *Overall Performance*: Fig. 6 compares the overall performance of different methods in KITTI sequence 00. HBST that performs feature match on raw image descriptors achieves the highest recall rate. However, both the computational cost and memory cost are too high to build the binary search tree so that it can be difficult to be implemented in practical application. In comparison, our approach provides a feasible solution for visual place recognition that is computationally efficient and accurate enough. Therefore, our approach achieves a good trade-off between both precision and efficiency. It is also notable that our method requires no offline training or vocabulary loading, making it more feasible for robotic real-time applications.

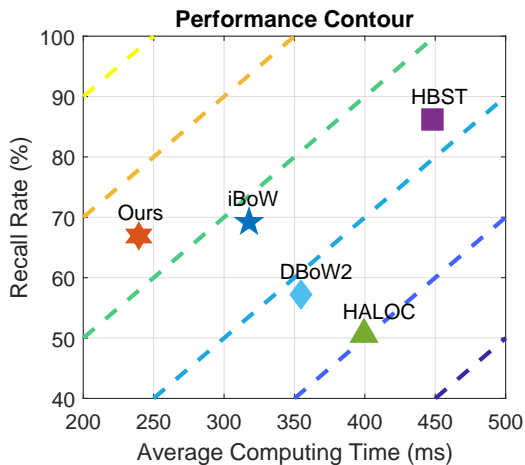


Fig. 6: Overall performance of different approaches.

V. CONCLUSION

In this paper, we present a novel framework for visual place recognition based on saliency re-identification. It mainly consists of two components, *i.e.* saliency detection and retrieval. To reduce the computational cost, both tasks are performed in frequency domain to take advantage of the efficiency of element-wise operation, resulting in an overall computational complexity of $\mathcal{O}(N \log N)$. The proposed method is open sourced as a C++ library. The experiments show that our method is more computationally efficient and accurate compared with the other state-of-the-art methods.

REFERENCES

- [1] 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, 2015.
- [2] H. Wang, C. Wang, and L. Xie, "Intensity scan context: Coding intensity and geometry relations for loop closure detection," in *IEEE international conference on robotics and automation (ICRA)*, 2020.
- [3] T. H. Nguyen, T.-M. Nguyen, M. Cao, and L. Xie, "Loosely-coupled ultra-wideband-aided scale correction for monocular visual odometry," *Unmanned Systems*, vol. 8, no. 02, pp. 179–190, 2020.
- [4] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [5] 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 *2015 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2015, pp. 6328–6335.
- [6] T. Cieslewski and D. Scaramuzza, "Efficient decentralized visual place recognition using a distributed inverted index," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 640–647, 2017.
- [7] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *Computer Vision (ICCV), 2011 IEEE international conference on*. IEEE, 2011, pp. 2564–2571.
- [8] S. Leutenegger, M. Chli, and R. Siegwart, "Brisk: Binary robust invariant scalable keypoints," in *2011 IEEE international conference on computer vision (ICCV)*. IEEE, 2011, pp. 2548–2555.
- [9] D. Schlegel and G. Grisetti, "Hbst: A hamming distance embedding binary search tree for visual place recognition," *arXiv preprint arXiv:1802.09261*, 2018.
- [10] E. Garcia-Fidalgo and A. Ortiz, "ibow-lcd: An appearance-based loop-closure detection approach using incremental bags of binary words," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3051–3057, 2018.
- [11] 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.
- [12] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [13] R. Cong, J. Lei, H. Fu, M.-M. Cheng, W. Lin, and Q. Huang, "Review of visual saliency detection with comprehensive information," *IEEE Transactions on Circuits and Systems for Video Technology*, 2018.
- [14] M. Donk and W. van Zoest, "Effects of salience are short-lived," *Psychological Science*, vol. 19, no. 7, pp. 733–739, 2008.
- [15] C. Wang, W. Wang, Y. Qiu, Y. Hu, and S. Scherer, "Visual Memorability for Robotic Interestingness via Unsupervised Online Learning," in *European Conference on Computer Vision (ECCV)*, 2020.
- [16] X. Hou and L. Zhang, "Saliency detection: A spectral residual approach," in *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*. IEEE, 2007, pp. 1–8.
- [17] 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 *Robotics and automation (ICRA), 2012 IEEE international conference on*. IEEE, 2012, pp. 4730–4735.
- [18] H. Bay, T. Tuytelaars, and L. Van Gool, "Surf: Speeded up robust features," in *European conference on computer vision*. Springer, 2006, pp. 404–417.
- [19] D. Gálvez-López and J. D. Tardos, "Bags of binary words for fast place recognition in image sequences," *IEEE Transactions on Robotics*, vol. 28, no. 5, pp. 1188–1197, 2012.
- [20] P. L. N. Carrasco, F. Bonin-Font, and G. Oliver-Codina, "Global image signature for visual loop-closure detection," *Autonomous Robots*, vol. 40, no. 8, pp. 1403–1417, 2016.
- [21] M. Gehrig, E. Stumm, T. Hinzmann, and R. Siegwart, "Visual place recognition with probabilistic voting," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3192–3199.
- [22] C. Wang, J. Yang, L. Xie, and J. Yuan, "Kervolutional neural networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 31–40.
- [23] Z. Chen, A. Jacobson, N. Sünderhauf, B. Upcroft, L. Liu, C. Shen, I. Reid, and M. Milford, "Deep learning features at scale for visual place recognition," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3223–3230.
- [24] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307.
- [25] S. Hausler, A. Jacobson, and M. Milford, "Multi-process fusion: Visual place recognition using multiple image processing methods," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 1924–1931, 2019.
- [26] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.
- [27] A. Srivastava, A. B. Lee, E. P. Simoncelli, and S.-C. Zhu, "On advances in statistical modeling of natural images," *Journal of mathematical imaging and vision*, vol. 18, no. 1, pp. 17–33, 2003.
- [28] B. S. Reddy and B. N. Chatterji, "An fft-based technique for translation, rotation, and scale-invariant image registration," *IEEE transactions on image processing*, vol. 5, no. 8, pp. 1266–1271, 1996.
- [29] C. Wang, L. Zhang, L. Xie, and J. Yuan, "Kernel cross-correlator," in *Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [30] C. Wang, T. Ji, T.-M. Nguyen, and L. Xie, "Correlation flow: robust optical flow using kernel cross-correlators," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 836–841.
- [31] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012.
- [32] 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.
- [33] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 year, 1000 km: The oxford robotcar dataset," *The International Journal of Robotics Research*, vol. 36, no. 1, pp. 3–15, 2017.
- [34] R. Mur-Artal and J. D. Tardós, "Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [35] X. Gao, R. Wang, N. Demmel, and D. Cremers, "Ldso: Direct sparse odometry with loop closure," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 2198–2204.