

Local positional graphs and attentive local features for a data and runtime-efficient hierarchical place recognition pipeline

Fangming Yuan, Stefan Schubert, Peter Protzel and Peer Neubert

Abstract—Large-scale applications of Visual Place Recognition (VPR) require computationally efficient approaches. Further, a well-balanced combination of data-based and training-free approaches can decrease the required amount of training data and effort and can reduce the influence of distribution shifts between the training and application phases. This paper proposes a runtime and data-efficient hierarchical VPR pipeline that extends existing approaches and presents novel ideas. There are three main contributions: First, we propose Local Positional Graphs (LPG), a training-free and runtime-efficient approach to encode spatial context information of local image features. LPG can be combined with existing local feature detectors and descriptors and considerably improves the image-matching quality compared to existing techniques in our experiments. Second, we present Attentive Local SPED (ATLAS), an extension of our previous local features approach with an attention module that improves the feature quality while maintaining high data efficiency. The influence of the proposed modifications is evaluated in an extensive ablation study. Third, we present a hierarchical pipeline that exploits hyperdimensional computing to use the same local features as holistic HDC-descriptors for fast candidate selection and for candidate reranking. We combine all contributions in a runtime and data-efficient VPR pipeline that shows benefits over the state-of-the-art method Patch-NetVLAD on a large collection of standard place recognition datasets with 15% better performance in VPR accuracy, 54× faster feature comparison speed, and 55× less descriptor storage occupancy, making our method promising for real-world high-performance large-scale VPR in changing environments. Code will be made available with publication of this paper.

I. INTRODUCTION

Visual place recognition (VPR) is a crucial component of SLAM systems. It finds the most similar images in a database with the given query image for various use cases such as loop closure detection or (re-)localization. The task of VPR is particularly challenging when the environmental condition between the database and the query set changes, e.g., from day to night or from summer to winter.

Convolutional neural networks (CNNs) based VPR methods can extract environmental condition- and viewpoint-robust image descriptors. These methods extract the descriptor of an image to either a single holistic (global) feature for the whole image [1] or a set of local features for regions of interest in the image [2]. Among CNN-based VPR methods, learning-based attentive local feature methods provide superior VPR performance. Compared to methods that densely extract large amounts of local features [3] [4], the attentive methods extract sparsely distributed and highly representative local features. However, existing best-performing learning-based attentive

This work was partially supported by the German Federal Ministry for Economic Affairs and Climate Action. F. Y., S. S. and P. P. are with the Chemnitz University of Technology, Germany. Email: {fangming.yuan, stefan.schubert, protzel}@ctf.tu-chemnitz.de. P. N. is with the University of Koblenz, Germany. Email: neubert@uni-koblenz.de.

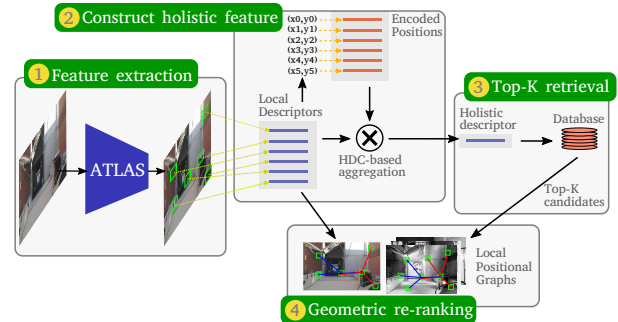


Fig. 1. An overview of the hierarchical VPR approach proposed in this paper: The pipeline first extracts local ATLAS image descriptors aggregated into a holistic image descriptor with hyperdimensional computing (HDC). These descriptors efficiently retrieve the top-K matching candidates before a final re-ranking with the local descriptors and the proposed Local Positional Graphs.

local feature methods require large-scaled [5] or expensively labeled [2] training datasets.

Geometric context among local features in an image significantly enhances the robustness and performance of VPR. Existing methods addressing lightweight image-wise geometric context [3] [6] are easy to implement but may have limited enhancement for the VPR performance. Others use learning-based methods to learn patch-wise or regional local feature geometric context representation [7][8], which are computationally expensive for inference.

Despite the advantages, local feature-based VPR is computationally expensive, hindering its usage for real-world large-scale VPR. Hierarchical VPR pipelines reduce the computational effort by retrieving the top-K candidates in the database with holistic features, then re-ranks the candidates by local features. However, the existing hierarchical VPR pipelines extract local features and holistic features with limited performance [3] or have lower query speed [9] [10] [3].

To address the above-mentioned limitations, this letter proposes the runtime and data-efficient hierarchical VPR pipeline illustrated in Fig. 1. First, inspired by the attention module of DELF [5], we extend our local feature approach LocalSPED-SoftMP [11] to the attentive local feature pipeline ATLAS. Compared to our previous work, ATLAS significantly increases the local feature performance in VPR with the same small (24K images) training dataset. Despite the 7× smaller training set, ATLAS shows a better pairwise mutual matching performance than DELF in our experiments. We conducted a detailed ablation study to address the differences between ATLAS and DELF. We found that the use of softmax normalization of the attention scores during training and the use of non-maximum suppression for local feature detection contribute significantly to ATLAS' high performance. For a second contribution, we propose a lightweight training-free algorithm

coined Local Positional Graph (LPG) to incorporate the local feature patch-wise geometric context for VPR. The proposed LPG algorithm shows advantages over several existing image-wise geometric context methods and significantly extends VPR performance for different local feature pipelines, especially for DELF. Third, we present a hierarchical VPR pipeline, that uses the same local features (like ATLAS or DELF) first in a hyperdimensional computing (HDC) based holistic descriptor [12] for fast candidate selection and then again for candidate re-ranking using the proposed LPG. Hir-ATLAS and Hir-Delf, the combinations with ATLAS and DELF, run up to $14\times$ faster than ATLAS and DELF but with only a 1.2% performance drop. We compare Hir-ATLAS and Hir-DELF with the hierarchical pipeline Patch-NetVLAD [3], demonstrating advantages in VPR accuracy (+15%), feature comparison speed ($54\times$), and feature storage occupancy ($55\times$).

II. RELATED WORK

Visual place recognition (VPR) is an active research topic. The recent paper [13] from 2023 provides a detailed introduction to VPR and an overview of relevant publications. Deep learning-based convolutional neural networks (CNNs) have contributed significantly to the performance improvement of VPR [14]. Especially for the creation of image descriptors, CNNs outperform handcrafted methods under severe appearance changes [15][16][2][5][3].

We can distinguish holistic (global) and local CNN descriptors. Both typically use the feature maps from an intermediate layer in a CNN as descriptors [17][16][15][18][19]. In earlier work, the regions of interest in an image for local feature-based methods were found either by handcrafted methods [15][18] or by detecting high activations in a pre-trained CNN [16][19]. However, the performance of these methods was limited because they were not specifically trained for VPR. To overcome this limitation, methods specifically trained for VPR began to emerge. In [1][20], a CNN-based holistic descriptor is trained for VPR with images collected from multiple webcams worldwide over a long period. To overcome the high memory occupation of the extracted local descriptors for the images [3][4], the authors propose attention methods to only detect and describe points of high interest [21][2][5]. In SuperPoint [21], the author train a pixel-level local feature pipeline on a synthetic dataset and warped real-world images using homographic adaptation to allow a self-supervised generation of known point correspondences between image pairs. To train D2-Net [2], the author uses structure from motion (SfM) to identify the ground truth point correspondences between images and subsequently use them to train a neural network for local feature extraction. Instead of using known point correspondences for training, DELF [5] uses a (weighted) linear combination of densely extracted local descriptors with a subsequent place classification to learn to detect and describe relevant local features. Based on [5] and [2], we proposed in our work [22][11] novel attentional pooling layers for a local feature detection and description pipeline. Similar to the training of DELF, we can avoid the creation of point correspondences between image pairs by treating the training as a place classification task with known

image correspondences. Our pipeline can achieve a high local descriptor performance using only a small training dataset.

Besides the local descriptor, the geometric context among the local features can also contribute to improve the performance of VPR [23][24][25][7][8]. However, the existing methods often use complex, time-consuming algorithms such as probabilistic inference [23] or random walk [26] for graph representation and comparison or introduce extra neural networks like a graph neural network [7] or a transformer [8] for geometric context feature extraction. However, there are also lightweight algorithms addressing the local feature geometric context. In [3], the author proposes the rapid spatial scoring method, which calculates the image similarity with the spatial displace of the mutually matched local features. In [6], the proposed Star-Hough algorithm constructs a star-shaped graph of local features in an image. However, both methods only address the image-wise geometric context among the local features, so they do not exploit patch-wise geometric details and are sensitive to viewpoint changes. In [4], the author only constructs the local spatial context between the dense extracted local image features in the nearby 3×3 -grid. Inspired by the above methods, we propose a novel lightweight graph method addressing patch-wise geometric context for image comparison.

Hierarchical VPR pipelines combine the advantages of holistic and local descriptors: Using the holistic descriptors to select the top-K candidates in the database for a query before re-ranking the candidates with the local features [9][3]. However, most holistic descriptor-based methods cannot extract local features [17][16][1][27][20][28]. In contrast, local descriptor methods do not provide holistic image representations [2][5]. Using different holistic and local features slows down the query speed. One versatile approach to aggregate local features in a global descriptor is hyperdimensional computing (HDC) [29]. HDC-DELF bundles a set of local descriptors with their position information into a single holistic descriptor, allowing the reuse of local descriptors for the holistic descriptors. In this letter, we use the local feature aggregation of [12] as an add-on module to generate holistic descriptors directly from our local descriptors to formulate a hierarchical VPR pipeline.

Despite their hierarchical structure, existing methods are often slow in the re-ranking stage [9] [3] [10] due to a large number of local features extracted per image. The proposed ATLAS pipeline extracts a condensed set of local features per image, significantly reducing the computing time for the re-ranking stage.

III. ALGORITHMIC DESCRIPTION

This section provides details of the VPR pipeline outlined in Fig. 1 and more detailedly illustrated in Fig. 2. We will first present the local feature extraction method ATLAS in Section III-A, followed by a description of the hierarchical HDC-based local-holistic VPR approach in Section III-B, and finally the details of the proposed Local Positional Graph (LPG) approach in Section III-C.

A. A detailed overview of ATLAS and its attention mechanism

The overall architecture of the proposed hierarchical pipeline is shown in Fig. 2. It comprises five main components:

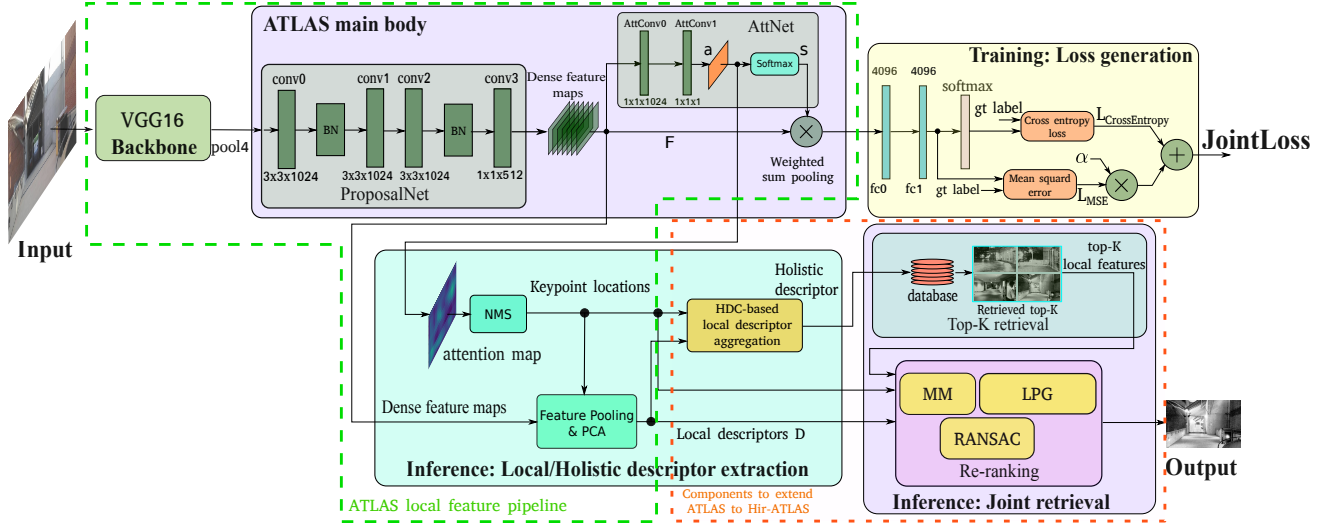


Fig. 2. An overview of the proposed hierarchical pipeline. The green dashed line covered area contains the ATLAS local feature pipeline components. In contrast, the orange dashed line covered area contains the components to extend ATLAS to Hir-ATLAS. Note that the top right yellow block is only used for ATLAS training, whereas the bottom part is exclusively used for inference, i.e., the actual application for VPR.

backbone, ATLAS main body, loss generation, local/holistic descriptor extraction, and joint retrieval. In the ATLAS local feature pipeline (green dashed line covered area), the backbone extracts the dense raw descriptors from the input image using the pool4-layer of VGG16 [30] trained on ImageNet [31]. Subsequently, these descriptors are further processed by the ATLAS main body, which is composed of the two trainable CNN networks *ProposalNet* and *AttNet* (Attention Net). *ProposalNet* transforms the raw descriptors from the backbone into the tensor $F \in \mathbb{R}^{H \times W \times C}$, which contains $H \times W$ VPR-specific dense local descriptors of dimensionality C . In ATLAS, we introduce *AttNet* into our pipeline, which is a two-layer CNN that extracts attention scores $a \in \mathbb{R}^{H \times W}$ from F to identify relevant, robust local features for VPR. The *softmax*-normalization is applied to all elements $a_{yx} \in a$ with

$$\forall y, x : s_{yx} = \frac{e^{a_{yx}}}{\sum_{i=1}^H \sum_{j=1}^W e^{a_{ij}}} \quad (1)$$

to obtain normalized attention scores $s \in \mathbb{R}^{H \times W}$. The final step of the ATLAS main body is the creation of a global representation $I \in \mathbb{R}^C$ of the input image using a weighted sum of all local descriptors in F :

$$\forall c : I_c = \sum_{y=1}^H \sum_{x=1}^W s_{yx} \cdot F_{yxc} \quad (2)$$

For **training**, I is passed to the *loss generation* module (cf. Fig. 2) to obtain the loss of the whole pipeline. Here, the training of ATLAS is formulated as a place classification task: I serves as input to two fully-connected layers (fc) that try to predict the correct place of the input image from N places. For the comparison of the predicted place and the ground truth, a combination of mean squared error L_{MSE} and cross entropy $L_{CrossEntropy}$ with

$$JointLoss = L_{CrossEntropy} + \alpha \cdot L_{MSE} \quad (3)$$

is used as proposed in [11]. The constant α weights the mean

squared error. Both losses are defined with

$$L_{CrossEntropy} = - \sum_{i=1}^N gt_i \cdot \log(\text{softmax}(v_i)) \quad (4)$$

$$L_{MSE} = \frac{1}{N} \sum_{i=1}^N (gt_i - v_i)^2 \quad (5)$$

The vectors $gt \in \mathbb{B}^N$ and $v \in \mathbb{R}^N$ represent the one-hot encoded ground truth place label of the input image and the predictions from the fully-connected layers (output of $fc1$ in Fig. 2).

During **inference**, the dense local descriptors F and the unnormalized attention scores a from the ATLAS main body are fed into the module for *local/holistic descriptor extraction* (cf. Fig. 2, bottom). Here, relevant local features with high attention scores are detected by a non-maximum suppression (NMS) with a 3×3 sliding window. To create the final set of local descriptors D , the descriptors in F are first pooled within a $d \times d$ -window around the found maxima in a before compressing the flattened patch features with principal component analysis (PCA) from dimensionality $d \cdot d \cdot C$ to a desired local descriptor length d_{loc} . Note that the PCA was previously learned on the training dataset. Section IV-C will provide an extensive ablation study to evaluate the design decisions.

B. Extending ATLAS for hierarchical VPR (Hir-ATLAS)

We extend ATLAS to Hir-ATLAS for hierarchical VPR by incorporating the local descriptor aggregation method of [12] into the module *local/holistic descriptor extraction* (cf. Fig. 2). The method uses hyperdimensional computing (HDC) to bundle the set D of local descriptors with their image positions into a single holistic descriptor (coined HDC-ATLAS). The *joint retrieval* module (cf. Fig. 2) can then perform the actual hierarchical VPR: In the first step, it uses the holistic descriptors to retrieve the top-K candidates from the database with the highest cosine similarity. The K candidates are re-ranked using the local descriptors in the second step. For re-ranking, we

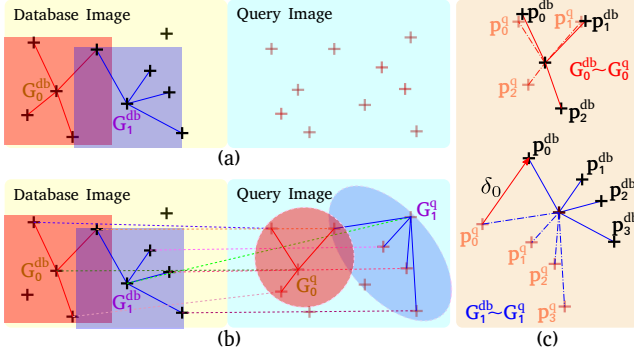


Fig. 3. Visualization of the Local Positional Graph (LPG). The crosses represent the local features in the images. (a) Creating a star-shaped graph for each local feature in the database image: Two graphs G_0^{db} and G_1^{db} are created for two local features. (b) Creation of corresponding graphs in the query image: The local features in the query image that are mutually matched to the root nodes in G_0^{db} and G_1^{db} serve as root nodes in the graphs G_0^q and G_1^q . Features in the query image that are mutually matched with the leaf nodes of G_0^{db} and G_1^{db} are leaf nodes in G_0^q and G_1^q . Unmatched leaves are discarded for G_0^{db} (bottom left node in G_0^{db}). (c) Graph comparison: The node positions in each graph are translated with the position of their root nodes to overlay the corresponding graphs. The displacement vectors δ_k are determined for all corresponding leaf nodes.

use three different approaches: 1) mutual matching (*MM*) of the local descriptors [15], 2) estimation of the fundamental matrix between two images using *RANSAC* [32], and 3) our proposed local positional graph (*LPG*) addressing the local feature patch-wise geometric context (Section III-C).

The three methods re-rank the image similarities $S_{db,q}$ between the query image and top-K candidates from the database with

$$S_{db,q} = \frac{1}{\sqrt{|D_{db}| \cdot |D_q|}} \sum_{\forall i,j} w_{ij} \cdot \cos(D_{db}^i, D_q^j) . \quad (6)$$

$D_{db} = \{D_{db}^i\}$ and $D_q = \{D_q^j\}$ are the sets of local descriptors from the db -th or q -th database or query image. $\cos(\cdot)$ is the cosine similarity between both input vectors. The weighting factor w_{ij} is set by the re-ranking method: For *MM*, $w_{ij} = 1$ if $\{D_{db}^i, D_q^j\}$ are mutual matches, otherwise $w_{ij} = 0$. For *RANSAC*, $w_{ij} = 1$ if the positions of the matched local features $\{i, j\}$ are inliers, otherwise $w_{ij} = 0$. How w_{ij} is set for *LPG* is explained in the following.

C. The local positional graph (LPG) for re-ranking

To overcome the limitations of *MM*, which does not consider the geometric context among local features, and of *RANSAC*, which is computationally very inefficient, we propose the local positional graph (*LPG*) for evaluating the similarity of two images. *LPG* is a lightweight training-free graph-based approach that exploits the patch-wise local feature geometric context to enhance the local feature performance in *VPR*. The *LPG* algorithm is composed of the three steps 1) graph creation, 2) graph comparison, and 3) image similarity evaluation.

1) *Graph creation*: For a database image with N local features, *LPG* first creates N star-shaped graphs G_n^{db} for each of the local features, where n indexes the graph constructed for local feature n . As shown in Fig. 3 (a), the local feature n in G_n^{db} is the root node, while all surrounding local features within

a rectangular window of size $h \times h$ (in local feature positional space) are leaf nodes. After a mutual matching of the local features between a database image and a query image, the corresponding star-shaped graphs G_n^q for the query image can be constructed directly based on G_n^{db} , as shown in Fig. 3 (b): The root node of G_n^q is the query image's local feature that mutually matches with the root node of G_n^{db} . The leaf nodes in G_n^q are the query image features that mutually match with the leaf nodes of G_n^{db} ; leaf nodes in G_n^{db} without a feature match in the query image are ignored in the successive graph comparison step. The outcome of the graph creation is a set of graph pairs $\{G_n^{db} \sim G_n^q\}$.

2) *Graph comparison*: For the graph comparison of a pair $\{G_n^{db} \sim G_n^q\}$, we first overlay their root nodes, as shown in Fig. 3 (c). Subsequently, a displacement vector δ_k is computed with

$$\delta_k = p_k^{db} - p_k^q \quad (7)$$

for each mutually matching leaf node pair with positions p_k^{db} and p_k^q . Here, $k \in [1..K]$ is the index of the K matched leaf nodes in the overlaid graph pair $\{G_n^{db} \sim G_n^q\}$. Next, the squared ℓ_2 -norm of each δ_k is mapped to displacement scores in the range $[0..1]$ with an unnormalized Gaussian

$$G(\delta_k) = \exp\left(-\frac{\|\delta_k\|_2^2}{2\sigma^2}\right) . \quad (8)$$

3) *Image similarity evaluation*: In the final step, we average $G(\delta_k)$ over all K matched leaf nodes of a graph pair $\{G_n^{db} \sim G_n^q\}$ with

$$w_{ij} = \frac{1}{K} \sum_{k=1}^K G(\delta_k) \quad (9)$$

to calculate the patch-wise geometric context. Here, i and j are the indices of the root node features in the database and query images, respectively. After the comparison of all graph pairs, all w_{ij} can be used to compute the final similarity $S_{db,q}$ of the image pair with Eq. (6). Note that $w_{ij} = 0$ if the i -th and j -th local features are not mutually matched.

The most time-consuming part of *LPG* is the graph creation for the database images. Fortunately, this can be done offline before the application for *VPR*. Moreover, the Gaussian function in Eq. (8) could be quantized to a reasonable resolution so that its calculation can be converted into a look-up table for efficient computation. The following experiments will demonstrate the low computational overhead of *LPG*.

IV. EXPERIMENTS

The experiments are conducted on a collection of standard visual place recognition datasets. First, we compare the local features of *ATLAS* and *DELf* and our predecessor local feature pipeline *LocalSPED-SoftMP (LSPD)*, using mutual matching (*MM*) and *LPG*. We also conduct the hyper-parameter investigation for *LPG* and compare *LPG* to three existing related algorithms. Second, we evaluate the holistic descriptor performance of *HDC-ATLAS*, *HDC-DELf* (holistic features aggregate with *DELf* local features), and the holistic feature used by *Patch-NetVLAD*. To compare with *Hir-ATLAS*, we extend *DELf* to *Hir-DELf* using the same methods described

TABLE I
MEAN AUC PERFORMANCES ON DIFFERENT DATASETS OF THE EXHAUSTIVE PAIRWISE IMAGE COMPARISON WITH LOCAL DESCRIPTORS (LEFT). COMPARISON OF THE LPG-RELATED METHOD (MIDDLE). ABLATION STUDY FOR ATLAS' LOCAL DESCRIPTOR ON MM (RIGHT). LSPD REPRESENTS THE METHOD LOCALSPED-SOFTMP [11]. RSS REPRESENTS THE RAPID SPATIAL SCORING OF [3]

Dataset	Local descriptor performance						Methods compared to LPG						Ablation study: ATLAS' local descriptor + MM						
	MM			LPG(ours)			RANSAC		Star-Hough [6]		POS [12]		RSS [3]		orig.	w/o softmax	w/o MSE	w/o NMS	w/o ProposalNet
	DELf[5]	LSPD[11]	ATLAS(ours)	DELf	LSPD	ATLAS	DELf	ATLAS	DELf	ATLAS	DELf	ATLAS	DELf	ATLAS					
GPW	0.58	0.75	0.85	0.87	0.85	0.92	0.81	0.90	0.84	0.89	0.83	0.84	0.65	0.79	0.85	0.08	0.82	0.76	0.76
Oxford	0.60	0.57	0.69	0.84	0.70	0.79	0.73	0.76	0.79	0.82	0.79	0.76	0.61	0.60	0.69	0.12	0.64	0.43	0.59
SFU	0.69	0.69	0.75	0.86	0.79	0.85	0.78	0.80	0.87	0.83	0.82	0.81	0.69	0.68	0.75	0.07	0.76	0.79	0.53
CMU	0.72	0.76	0.77	0.80	0.78	0.79	0.77	0.78	0.79	0.77	0.77	0.77	0.55	0.69	0.77	0.10	0.77	0.63	0.71
Nordland	0.50	0.68	0.75	0.78	0.80	0.87	0.68	0.81	0.76	0.84	0.73	0.82	0.51	0.70	0.75	0.08	0.77	0.52	0.57
StLucia	0.47	0.47	0.50	0.69	0.55	0.62	0.54	0.53	0.67	0.55	0.59	0.55	0.50	0.44	0.50	0.06	0.50	0.33	0.43
mean	0.60	0.65	0.72	0.81	0.74	0.81	0.72	0.76	0.79	0.79	0.76	0.76	0.58	0.65	0.72	0.08	0.70	0.56	0.60

TABLE II
AUC OF ATLAS AND DELF LOCAL FEATURES IN GPW DAY-LEFT VS. NIGHT-RIGHT SEQUENCE WITH DIFFERENT LPG HYPER-PARAMETER COMBINATIONS OF $h \times h$ AND σ . BEST COMBINATION FOR BOTH ATLAS AND DELF IS BOLD.

h×h	10×10		20×20		40×40		60×60		80×80	
	ATLAS	DELf	ATLAS	DELf	ATLAS	DELf	ATLAS	DELf	ATLAS	DELf
$\sigma = 0.5$	0.77	0.35	0.85	0.55	0.87	0.66	0.88	0.74	0.88	0.77
$\sigma = 1.0$	0.83	0.46	0.87	0.66	0.88	0.75	0.88	0.78	0.87	0.78
$\sigma = 2.0$	0.83	0.50	0.86	0.66	0.87	0.73	0.86	0.75	0.85	0.74
$\sigma = 3.0$	0.83	0.48	0.85	0.62	0.85	0.69	0.85	0.72	0.83	0.71

in Section III-B. Next, we evaluate the hierarchical VPR performance of Hir-ATLAS, Hir-DELf, and Patch-NetVLAD on three aspects: VPR performance, query speed, and disk storage occupancy. Finally, we conduct an ablation study for the ATLAS local feature pipeline.

A. Experimental setup

1) *Parameters*: We use all ATLAS local features with their local descriptors. The local features are detected with the non-maximum suppression (NMS) on the attention scores a using a 3×3 sliding window. To create the local descriptors D^i , we pool the dense local features in F in a 7×7 -window ($d = 7$) around a local maximum in a before they are compressed with PCA to dimensionality $d_{loc} = 1024$. For DELf local features, we use the implementation provided by the authors and select the top 200 local descriptors (with the highest scores) of dimensionality 1024. For LPG, we normalize the feature positions to $[0..100)$. We also quantize $\|\delta_k\|_2^2$ in this local feature coordinate space to sample the Gaussian function in Eq. (8). We conduct a hyper-parameter search on *day-left - night-right* sequence of the GPW dataset to find the best σ and h combination for both ATLAS and DELf. As shown in Table II, the LPG parameters are set to $\sigma = 1.0$ and $h = 60$ for maximum performance. For Patch-NetVLAD, we also use the author's model, which is trained on the mapillary dataset. For 4096-dimensional HDC-ATLAS holistic descriptors, we use the HDC-based local descriptor aggregation method from [12] with $n_x = 5$ and $n_y = 9$. We use the holistic descriptor HDC-DELf from [12] without the feature normalization for the holistic feature of Hir-DELf.

2) *Training*: We train ATLAS with the same training procedure and dataset as in our previous work [11].

3) *Benchmark*: We evaluate all methods on the six different datasets GardensPoint Walking (GPW) [33], StLucia (multiple times of day) [34], Oxford RobotCar [35], CMU [36], Nordland [37], and SFU Mountain [38]. StLucia, Oxford, and CMU were sampled with one frame per second. All images of each sequence in GPW and SFU are used. We sampled 1,000 images of unique places (without tunnels) for Nordland as in [12].

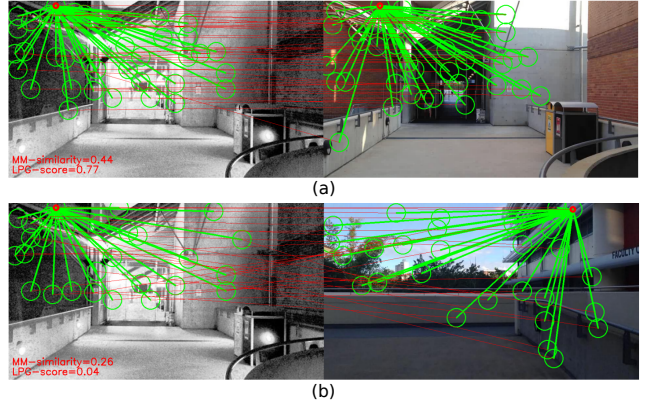


Fig. 4. LPG graphs with correctly matched root nodes (a) and incorrectly matched root nodes (b) using local ATLAS features between database images (left) and query images (right). Red circles show the root nodes. Green circles show leaf nodes with a radius that corresponds to the pixel size in the attention map. Red lines connect matched leaf nodes.

TABLE III
MEAN AUC/RECALL@100 ($K = 100$) PERFORMANCES ON DIFFERENT DATASETS OF THE HOLISTIC DESCRIPTORS HDC-ATLAS FOR HIR-ATLAS, HDC-DELf FOR HIR-DELf, AND HOL-PVLAD FOR PATCH-NETVLAD.

Dataset	HDC-ATLAS (ours)		HDC-DELf [12]		HOL-PVLAD [3]	
	AUC	R@100	AUC	R@100	AUC	R@100
GPW	0.52	0.98	0.63	1.00	0.44	0.97
Oxford	0.35	0.94	0.64	0.95	0.55	0.90
SFU	0.58	0.98	0.61	0.98	0.06	0.90
CMU	0.65	0.89	0.68	0.89	0.51	0.90
Nordland	0.64	0.99	0.61	0.98	0.15	0.88
StLucia	0.41	0.83	0.42	0.87	0.08	0.63
mean	0.51	0.94	0.61	0.94	0.33	0.87

4) *Metrics*: For performance evaluation, we use the area under the precision-recall curve (AUC) and Recall@K [13].

B. Local feature comparison and LPG evaluation

We first evaluate ATLAS' local descriptor's performance for VPR and compare it with LSPD and DELf. Therefore, we use the local descriptors to compare all query images with the whole database exhaustively. We use mutual matching (MM) and LPG to calculate the similarity of each image pair. In addition, we compare LPG with three image-wise geometric context-based methods: Star-Hough [6], POS (as used in [12]), and the Rapid Spatial Scoring (RSS) as used in Patch-NetVLAD [3].

Table I (left) shows the obtained mean AUC performances on different datasets of the three local feature pipelines with MM and LPG. For each dataset, we use the same sequences as in Table IV. ATLAS with MM outperforms both LSPD and DELf on most of the datasets. This is also reflected in the mean, best- and worst-case performances over all datasets: ATLAS' local descriptor with MM outperforms LSPD by almost 11% and

TABLE IV

AUC PERFORMANCES OF THE HIERARCHICAL VPR PIPELINES. HIR-ATLAS WITH NVD+RANSAC USES PATCH-NETVLAD’S HOLISTIC DESCRIPTOR FOR CANDIDATE SELECTION BEFORE RE-RANKING WITH THE ATLAS’ LOCAL DESCRIPTORS AND RANSAC.

Dataset	Query - DB	Hir-DELFL top-100			Hir-ATLAS top-100 (ours)				Patch-NetVLAD top-100	
		MM	LPG (ours)	RANSAC	MM	LPG (ours)	RANSAC	NV+RANSAC	Performance	Speed
GPW	day-left-day-right	0.92	0.98	0.98	0.97	0.98	0.98	0.99	1.00	1.00
	day-right-night-right (night)	0.49	0.86	0.79	0.79	0.89	0.85	0.84	0.92	0.74
	day-left-night-right (night)	0.26	0.78	0.66	0.75	0.85	0.84	0.83	0.80	0.58
Oxford	14-12-09-15-05-19	0.91	0.95	0.94	0.70	0.79	0.77	0.82	0.67	0.61
	14-12-09-15-08-28	0.34	0.75	0.56	0.48	0.62	0.56	0.62	0.67	0.47
	14-12-09-14-11-25	0.77	0.89	0.82	0.88	0.88	0.89	0.91	0.65	0.60
	14-12-09-14-12-16 (night)	0.30	0.66	0.57	0.68	0.70	0.70	0.67	0.44	0.28
	15-05-19-15-02-03	0.81	0.96	0.94	0.68	0.93	0.92	0.92	0.38	0.32
	15-08-28-14-11-25	0.59	0.79	0.71	0.86	0.72	0.64	0.67	0.51	0.38
SFU	dry-dusk	0.76	0.86	0.80	0.73	0.82	0.78	0.78	0.93	0.86
	dry-jan	0.58	0.85	0.76	0.77	0.86	0.82	0.78	0.86	0.72
	dry-wet	0.72	0.86	0.78	0.75	0.85	0.81	0.78	0.92	0.82
CMU	20110421-20100901	0.79	0.83	0.80	0.81	0.83	0.82	0.83	0.65	0.54
	20110421-20100915	0.74	0.78	0.77	0.78	0.79	0.78	0.79	0.55	0.47
	20110421-20101221	0.63	0.69	0.68	0.65	0.68	0.66	0.76	0.56	0.46
	20110421-20110202	0.72	0.85	0.80	0.83	0.85	0.83	0.83	0.61	0.57
Nordland	spring-winter	0.54	0.92	0.79	0.79	0.94	0.86	0.79	0.85	0.80
	spring-summer	0.45	0.77	0.70	0.80	0.88	0.84	0.83	0.92	0.87
	summer-winter	0.17	0.53	0.40	0.54	0.74	0.63	0.57	0.68	0.55
	summer-fall	0.83	0.93	0.92	0.93	0.94	0.94	0.93	0.97	0.96
StLucia	100909-0845-180809-1545	0.34	0.54	0.41	0.38	0.50	0.41	0.25	0.30	0.24
	100909-1000-190809-1410	0.53	0.73	0.60	0.54	0.65	0.57	0.47	0.50	0.45
	100909-1210-210809-1210	0.69	0.76	0.70	0.64	0.68	0.65	0.68	0.75	0.72
best		0.92	0.98	0.98	0.97	0.98	0.98	0.99	1.00	1.00
worst		0.17	0.53	0.40	0.38	0.50	0.41	0.25	0.30	0.24
mean		0.60	0.81	0.73	0.73	0.80	0.76	0.75	0.70	0.61

DELFL by even 20% in mean AUC. The LPG algorithm further enhances the local feature performance over MM for all three descriptors: The mean AUC of ATLAS, LSPD, and DELFL is improved by approximately 13%, 14% and 35%. Interestingly, with LPG, DELFL reaches the same performance as ATLAS. One possible explanation is that ATLAS already learned to incorporate more geometric context during training than DELFL, which then leaves less room for further improvement for ATLAS. On the other hand, DELFL can more strongly benefit from the additional geometric context provided by LPG.

The middle of table Table I shows the evaluation results of the three local feature image-wise geometric context-based methods Star-Hough, POS, and RSS. The numbers are similar for ATLAS and DELFL. Only with RSS, there is a noticeable difference in favour of ATLAS. The results show that LPG provides better performance than the other three context-based methods.

Fig. 4 shows two pairs of LPG graphs between a database image and a query image. The root nodes are actual mutual matches that were matched either correctly (true positive) or incorrectly (false positive). Both examples demonstrate how LPG contributes to a better performance: In the correct example (Fig. 4a), LPG further increases the similarity between the root nodes using the leaf nodes, while in the incorrect example (Fig. 4b), LPG further decreases the similarity of the root nodes using the leaf nodes.

C. Ablation study of ATLAS’ local descriptor

In this section, we conduct an ablation study to find out the critical components of the local descriptor pipeline ATLAS, which are different from the local descriptor DELFL [5] and contribute to the high performance of ATLAS. The four key components that only appear in ATLAS but not in DELFL are (cf. Fig. 2)

- 1) a *softmax* normalization for the attention scores from the attention network AttNet (cf. Eq. (1)),

- 2) the mean squared error L_{MSE} (*MSE*) that is used in addition to the cross entropy loss during training (cf. Eq. (3)),
- 3) a non-maximum suppression (*NMS*) of the attention scores to select relevant local features instead of using the K highest scored local features as in DELFL, and
- 4) *ProposalNet* that post-processes the dense local features from VGG16.

For the ablation study, we run four experiments and omit either component 1, 2, or 4 or replace component 3 with a selection of the 200 highest scored local features as in DELFL. We omit 4 by extracting the local feature from VGG16 pool4 layer output feature maps (the input feature maps for ProposalNet) instead of ProposalNet’s output according to the detected keypoint locations. The modified pipelines addressing 1 and 2 are retrained with the same training procedure as with ATLAS (cf. Section IV-A2). We then evaluate the local descriptors with the trained modified pipeline and perform an exhaustive pairwise image comparison between the query set and the database using MM as described in Section IV-A. Table I (right) shows the obtained results with the original local ATLAS descriptor (column orig.) and with the four modified versions. The performance measurements show that the softmax normalization (component 1) is essential during training. Without it, the pipeline fails to learn a meaningful local descriptor for VPR, so the performance drops to nearly zero. The mean squared error (component 2) barely contributes to the performance and only yields a 3% performance gain. However, we observed in the experiment that the additional mean squared error term could increase the sparsity of the attention scores after NMS: When trained with the combined loss function (i.e., $\alpha = 10$ in Eq. (3)), the pipeline extracts 15% less local descriptors but provides better VPR performance. Omitting the non-maximum suppression (component 3) or ProposalNet (component 4) also clearly decreases the performance of the local ATLAS descriptors by 22% or 17%, respectively, indicating both contribute to ATLAS. The ablation study

indicates, the attention score softmax normalization and the joint loss function for training, and the use of NMS for local feature detection jointly contribute to the performance of ATLAS.

D. Evaluations of Hir-ATLAS, Hir-DELF, and Patch-NetVLAD

In this section, we evaluate the performance of Hir-ATLAS for hierarchical VPR as described in Section III and compare it with Hir-DELF and the state-of-the-art hierarchical approach Patch-NetVLAD [3]. All three methods retrieve and re-rank the top 100 image pairs per query. Patch-NetVLAD is evaluating either in *performance mode* or *speed mode*: Both use NetVLAD [28] based holistic feature (here we name this holistic feature HOL-PVLAD) to select the top-100 candidates. For re-ranking, the *performance mode* builds upon the multi-scale image patch feature and uses RANSAC-based fundamental matrix estimation for outlier rejection, while the *speed mode* only uses a single-scale patch feature and the RSS algorithm. We use MM, LPG or RANSAC (cf. Sections III-B and III-C) to re-rank either Hir-ATLAS or Hir-DELF.

In the first experiment, we measure the candidate selection performance of the holistic descriptors HDC-ATLAS, HDC-DELF, and HOL-PVLAD used for the three methods Hir-ATLAS, Hir-DELF, and Patch-NetVLAD. Table III shows the corresponding AUC and Recall@100 ($K = 100$) values. On average, HDC-DELF slightly outperforms HDC-ATLAS on both metrics, while HOL-PVLAD performs worse than HDC-DELF and HDC-ATLAS.

Next, we evaluate the VPR performance of the three hierarchical methods. The obtained results are shown in Table IV. On average, Hir-ATLAS performs best, outperforms Hir-DELF on MM and RANSAC configurations, and outperforms Patch-NetVLAD in both *speed* and *performance* mode. Although Hir-DELF uses the best-performing holistic descriptor HDC-DELF for candidate selection, Hir-DELF in LPG configuration takes only a tiny advantage in mean AUC than Hir-ATLAS with LPG. Compared to Patch-NetVLAD, Hir-ATLAS with LPG achieves approx. 14% to 31% higher mean AUC performance. Only on the SFU dataset, Patch-NetVLAD is best suited and outperforms all other methods. We also evaluate Hir-ATLAS using the holistic descriptor HOL-PVLAD for top-100 candidate selection and using ATLAS’ local feature with RANSAC for re-ranking (referred to as NV+RANSAC) to allow a comparison of the local descriptors between ATLAS and Patch-NetVLAD. In this setup, Hir-ATLAS (with NV+RANSAC) still performs 7% better in mean AUC than Patch-NetVLAD in *performance* setup, which implies a better local feature performance of Hir-ATLAS than Patch-NetVLAD for VPR.

We finally compare the quality of the re-ranking between the best-performing pipeline setups of Hir-ATLAS, Hir-DELF, and Patch-NetVLAD. We, therefore, measure the Recall@K performance with $K = \{1, 5, 10, 20\}$ after re-ranking the top 100 candidates. Table V shows the obtained average recall of the entire benchmark datasets: On average, Hir-DELF with LPG performs best for all K values. Hir-DELF and Hir-ATLAS share similar Recall@K performance for most of the k values and outperform Patch-NetVLAD in *performance mode*.

In summary, our pipeline Hir-ATLAS outperforms Patch-NetVLAD in all re-ranking methods and outperforms Hir-DELF

TABLE V
AVERAGE RECALL@K OVER ALL BENCHMARK DATASETS.

Pipeline	R@1	R@5	R@10	R@20
Hir-ATLAS top-100 LPG (ours)	0.55	0.68	0.74	0.80
Hir-ATLAS top-100 RANSAC (ours)	0.53	0.67	0.73	0.80
Hir-DELF top-100 LPG	0.54	0.69	0.74	0.81
Hir-DELF top-100 RANSAC	0.50	0.66	0.72	0.80
Patch-NetVLAD top-100 performance	0.52	0.63	0.68	0.74

in MM and RANSAC. The LPG significantly improves the performance of Hir-DELF and reaches the same performance as Hir-ATLAS in hierarchical VPR.

E. Runtime and disk usage

Besides the AUC value and recall, other essential factors for practical applications are the query speed and disk space to store the image features. The feature comparison between query image features and the database image features takes the majority of query time. Therefore, we measured only the feature comparison time for the query. We measure the feature comparison time on the entire benchmark with up to 27K queries. We use a computation platform with Intel i7-4790 CPU and NVIDIA RTX4090 GPU. Table VI shows the obtained total feature comparison time, average feature comparison latency, relative speed-up reference to PatchNetVLAD in performance mode, feature disk space occupancy to store all the benchmark image features, and the corresponding mean AUC (mAUC) value of the three pipelines with different setups. Patch-NetVLAD in *performance* mode takes 1,400GB disk space to store all the benchmark image features, which is $55\times$ larger than Hir-ATLAS and Hir-DELF both use roughly 25.5GB disk space. This is because Patch-NetVLAD in performance mode extracts 2,826 4096-dimensional local features per image, while ATLAS and DELF extract, on average, only 200 1024-dimensional local features per image.

Patch-NetVLAD in *speed* mode provides lower VPR performance, requires more memory, and is not as fast as Hir-ATLAS and Hir-DELF in LPG mode. Compared to Hir-ATLAS and Hir-DELF in RANSAC setup, Patch-NetVLAD in performance mode runs $8\times$ and $4.68\times$ slower. The best-performing Hir-ATLAS and Hir-DELF setup (Top-100 LPG) runs even $54\times$ and $74\times$ faster than the best-performing Patch-NetVLAD setup (top-100 performance). The average feature comparison latency of setup Top-100 LPG for Hir-ATLAS and Hir-DELF takes less than 27ms, which can be regarded as real-time for many practical applications. For the comparison between three re-ranking methods MM, LPG, and RANSAC in Hir-ATLAS and Hir-DELF, LPG is $1.4\times$ slower than MM but provides better performance. LPG also performs better than RANSAC and runs $6.7\times$ and $15.8\times$ faster for Hir-ATLAS and Hir-DELF. Finally, the results demonstrate that the hierarchical approach is faster than the exhaustive approaches: The hierarchical approach with top-100 candidates runs approximately $14\times$ faster than the exhaustive local descriptor comparison but with negligible performance drops. Despite different use cases, we also compared ATLAS+LPG and LightGlue [7] with local SuperPoint features [21] on the GPW and SFU datasets and found that ATLAS runs $150\times$ to $400\times$ faster than LightGlue, is more robust to viewpoint changes, and requires significantly fewer local features.

TABLE VI
FEATURE COMPARISON TIME AND DISK USAGE OF DIFFERENT PIPELINES WITH DIFFERENT SETUPS. THE SPEED UP COLUMN USES PATCH-NETVLAD TOP-100 PERFORMANCE AS A REFERENCE.

Setup	Total comparison time	Avg latency	Speed up	Disk usage	mAUC
Patch-NetVLAD top-100 performance	~11.1h	1400ms	1×	~1400GB	0.70
Hir-ATLAS top-100 RANSAC (ours)	4956s	180ms	8×	~25.5GB	0.76
Hir-DELf top-100 RANSAC	8538s	310ms	4.68×	~25.5GB	0.73
Patch-NetVLAD top-100 speed	2594s	95ms	14.7×	~419GB	0.60
Hir-ATLAS top-100 LPG (ours)	741s	27ms	54×	~25.5GB	0.80
Hir-DELf top-100 LPG	542s	20ms	74×	~25.5GB	0.81
Hir-ATLAS top-100 MM (ours)	526s	19ms	75×	~25.5GB	0.73
Hir-DELf top-100 MM	388s	14ms	103×	~25.5GB	0.60
ATLAS exhaustive MM (ours)	7383s	270ms	5.4×	~25.5GB	0.72
DELf exhaustive MM	5170s	190ms	7.7×	~25.5GB	0.60

V. CONCLUSION

In this letter, we propose the local feature pipeline ATLAS which achieves a decent VPR performance based on a small training dataset. The proposed LPG algorithm demonstrates a significant VPR performance enhancement for state-of-the-art local feature pipelines over existing spatial context-based methods. Using LPG and the HDC method we extend ATLAS to Hir-ATLAS and DELf to Hir-DELf for hierarchical VPR. The high performance, low training cost, low latency for query, and small disk occupancy make Hir-ATLAS and Hir-DELf promising for real-world large-scale VPR in changing environments.

REFERENCES

- [1] 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 *International Conference on Robotics and Automation (ICRA)*, 2017, pp. 3223–3230.
- [2] M. Dusmanu, I. Rocco, T. Pajdla, M. Pollefeys, J. Sivic, A. Torii, and T. Sattler, "D2-Net: A trainable CNN for joint description and detection of local features," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 8084–8093.
- [3] S. Hausler, S. Garg, M. Xu, M. Milford, and T. Fischer, "Patch-NetVLAD: Multi-scale fusion of locally-global descriptors for place recognition," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 14 141–14 152.
- [4] L. G. Camara and L. Přeučil, "Spatio-semantic ConvNet-based visual place recognition," in *European Conference on Mobile Robots (ECMR)*, 2019, pp. 1–8.
- [5] H. Noh, A. Araujo, J. Sim, T. Weyand, and B. Han, "Large-scale image retrieval with attentive deep local features," in *International Conference on Computer Vision (ICCV)*, 2017, pp. 3476–3485.
- [6] P. Neubert and P. Protzel, "Local region detector + CNN based landmarks for practical place recognition in changing environments," in *European Conference on Mobile Robots (ECMR)*, 2015, pp. 1–6.
- [7] P. Lindenberger, P.-E. Sarlin, and M. Pollefeys, "LightGlue: Local feature matching at light speed," in *International Conference on Computer Vision (ICCV)*, 2023, pp. 17 627–17 638.
- [8] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, "LoFTR: Detector-free local feature matching with transformers," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8922–8931.
- [9] B. Cao, A. Araujo, and J. Sim, "Unifying deep local and global features for image search," in *European Conf. on Computer Vision (ECCV)*, 2020.
- [10] R. Wang, Y. Shen, W. Zuo, S. Zhou, and N. Zheng, "TransVPR: Transformer-based place recognition with multi-level attention aggregation," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 13 638–13 647.
- [11] F. Yuan, P. Neubert, S. Schubert, and P. Protzel, "SoftMP: Attentive feature pooling for joint local feature detection and description for place recognition in changing environments," in *International Conference on Robotics and Automation (ICRA)*, 2021, pp. 5847–5853.
- [12] P. Neubert and S. Schubert, "Hyperdimensional computing as a framework for systematic aggregation of image descriptors," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [13] S. Schubert, P. Neubert, S. Garg, M. Milford, and T. Fischer, "Visual place recognition: A tutorial," *IEEE Robotics and Automation Magazine (RAM) (to appear)*, 2023.
- [14] C. Masone and B. Caputo, "A survey on deep visual place recognition," *IEEE Access*, vol. 9, pp. 19 516–19 547, 2021.
- [15] 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," in *Robotics: Science and Systems (RSS)*, 2015.
- [16] Z. Chen, F. Maffra, I. Sa, and M. Chli, "Only look once, mining distinctive landmarks from ConvNet for visual place recognition," in *International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 9–16.
- [17] N. Sünderhauf, S. Shirazi, F. Dayoub, B. Upcroft, and M. Milford, "On the performance of ConvNet features for place recognition," in *International Conference on Intelligent Robots and Systems (IROS)*, 2015.
- [18] P. Neubert, "Superpixels and their application for visual place recognition in changing environments," PhD Thesis, Chemnitz University of Technology, 2015. [Online]. Available: <http://nbn-resolving.de/urn:nbn:de:bsz:ch1-qucosa-190241>
- [19] A. Khaliq, S. Ehsan, Z. Chen, M. Milford, and K. McDonald-Maier, "A holistic visual place recognition approach using lightweight CNNs for significant viewpoint and appearance changes," *IEEE Transactions on Robotics*, vol. 36, no. 2, pp. 561–569, 2020.
- [20] Z. Chen, L. Liu, I. Sa, Z. Ge, and M. Chli, "Learning context flexible attention model for long-term visual place recognition," *IEEE Robotics and Automation Letters (RA-L)*, vol. 3, no. 4, pp. 4015–4022, 2018.
- [21] D. DeTone, T. Malisiewicz, and A. Rabinovich, "SuperPoint: Self-supervised interest point detection and description," in *Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018, pp. 337–349.
- [22] F. Yuan, P. Neubert, and P. Protzel, "LocalSPED: A classification pipeline that can learn local features for place recognition using a small training set," in *Towards Autonomous Robotic Systems (TAROS)*, 2020.
- [23] E. Stumm, C. Mei, S. Lacroix, J. Nieto, M. Hutter, and R. Siegwart, "Robust visual place recognition with graph kernels," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4535–4544.
- [24] P. Gao and H. Zhang, "Long-term place recognition through worst-case graph matching to integrate landmark appearances and spatial relationships," in *International Conference on Robotics and Automation (ICRA)*, 2020, pp. 1070–1076.
- [25] Z. Luo, T. Shen, L. Zhou, J. Zhang, Y. Yao, S. Li, T. Fang, and L. Quan, "ContextDesc: Local descriptor augmentation with cross-modality context," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2522–2531.
- [26] A. Gawel, C. D. Don, R. Siegwart, J. Nieto, and C. Cadena, "X-View: Graph-based semantic multi-view localization," *IEEE Robotics and Automation Letters (RA-L)*, vol. 3, no. 3, pp. 1687–1694, 2018.
- [27] A. Ali-bey, B. Chaib-draa, and P. Giguère, "MixVPR: Feature mixing for visual place recognition," in *Winter Conference on Applications of Computer Vision*, 2023, pp. 2998–3007.
- [28] R. Arandjelović, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 6, pp. 1437–1451, 2018.
- [29] P. Neubert, S. Schubert, and P. Protzel, "An introduction to hyperdimensional computing for robotics," *KI - Künstliche Intelligenz*, vol. 33, no. 4, pp. 319–330, 2019.
- [30] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations (ICLR)*, 2015.
- [31] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.
- [32] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge University Press, 2004.
- [33] A. Glover, "Day and night with lateral pose change datasets," 2014.
- [34] A. J. Glover, W. P. Maddern, M. J. Milford, and G. F. Wyeth, "FAB-MAP + RatSLAM: Appearance-based SLAM for multiple times of day," in *International Conference on Robotics and Automation (ICRA)*, 2010.
- [35] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 year, 1000km: The oxford RobotCar dataset," *The International Journal of Robotics Research*, vol. 36, no. 1, pp. 3–15, 2017.
- [36] H. Badino, D. Huber, and T. Kanade, "Visual topometric localization," in *Intelligent Vehicles Symposium (IV)*, 2011, pp. 794–799.
- [37] N. Sünderhauf, P. Neubert, and P. Protzel, "Are we there yet? Challenging SeqSLAM on a 3000 km journey across all four seasons," in *International Conference on Robotics and Automation (ICRA) Workshop on Long-Term Autonomy*, 2013, pp. 1–3.
- [38] J. Bruce, J. Wawerla, and R. Vaughan, "The SFU mountain dataset: Semi-structured woodland trails under changing environmental conditions," in *International Conference on Robotics and Automation (ICRA) Workshop on Visual Place Recognition in Changing Environments*, 2015.