

Visual Place Recognition for Aerial Robotics: Exploring Accuracy-Computation Trade-off for Local Image Descriptors

1st Bruno Ferrarini
CSEE School
University of Essex
Colchester, CO4 3SQ, UK
bferra@essex.ac.uk

2nd Maria Waheed
COMSATS University
Islamabad, Pakistan
maria.waheed97@gmail.com

3rd Sania Waheed
National University of
Sciences and Technology (NUST)
Islamabad, Pakistan
saniawaheed97@gmail.com

4th Shoaib Ehsan
CSEE School
University of Essex
Colchester, CO4 3SQ, UK
sehsan@essex.ac.uk

5th Michael Milford
Science and Engineering Faculty
Queensland University of Technology
Brisbane, Australia
michael.milford@qut.edu.au

6th Klaus D. McDonald-Maier
CSEE School
University of Essex
Colchester, CO4 3SQ, UK
kdm@essex.ac.uk

Abstract—Visual Place Recognition (VPR) is a fundamental yet challenging task for small Unmanned Aerial Vehicle (UAV). The core reasons are the extreme viewpoint changes, and limited computational power onboard a UAV which restricts the applicability of robust but computation intensive state-of-the-art VPR methods. In this context, a viable approach is to use local image descriptors for performing VPR as these can be computed relatively efficiently without the need of any special hardware, such as a GPU. However, the choice of a local feature descriptor is not trivial and calls for a detailed investigation as there is a trade-off between VPR accuracy and the required computational effort. To fill this research gap, this paper examines the performance of several state-of-the-art local feature descriptors, both from accuracy and computational perspectives, specifically for VPR application utilizing standard aerial datasets. The presented results confirm that a trade-off between accuracy and computational effort is inevitable while executing VPR on resource-constrained hardware.

Index Terms—Local Image Descriptors, Visual Place Recognition, Comparison, Unmanned Aerial Vehicles

I. INTRODUCTION

Autonomous navigation of UAVs has been receiving great attention recently [1]–[3] as it has a wide variety of industrial applications, such as aerial imaging and surveying [4]. As part of autonomous navigation, place recognition is critical for UAV localization [5]. When the estimation of a vehicle position drifts due to accumulated errors over time, re-localization is possible when a reference location/already-visited place is detected [6]. Place recognition for a UAV is usually addressed using visual information, hence called visual place recognition. This is motivated by the availability, cost, size and weight of modern cameras, which make them feasible to be installed even on a small aerial vehicle.

This work has been supported by the UK Engineering and Physical Sciences Research Council EPSRC [EP/K004638/1, EP/R02572X/1 and EP/P017487/1]

VPR is a particularly challenging problem for a small UAV due to the extreme viewpoint changes and limited computational power onboard [2]. State-of-the-art VPR methods are only robust to small viewpoint changes [7], [8] or their computation intensive nature [9] makes their use prohibitive for small UAVs. Some recent works [2], [10] proposed VPR pipelines based on local image features as they can be extracted efficiently on resource-constrained hardware. As noted in [10], the choice of the local feature descriptor implies a trade-off between the computation efficiency and the accuracy of the overall visual place recognition system. This calls for detailed investigation into accuracy-computation trade-off for local feature descriptors specifically for VPR application.

To this end, this paper explores the accuracy-computation trade-off of several state-of-the-art local feature descriptors for VPR under mild to extreme viewpoint changes in small UAVs using standard ground-aerial image datasets [2]. Precision-Recall and computation time are used as the basis of exploring this accuracy-computation trade-off. Results are presented for SIFT [11], SURF [12], BRISK [13], AKAZE [14] and ORB [15] descriptors. However, the proposed evaluation method can be extended to any local image feature descriptor. The results presented confirm that the more computationally expensive descriptors yield to better accuracy at the cost of a longer localization time. A trade-off between accuracy and computation effort appears inevitable while executing VPR on a limited resource hardware.

The rest of this paper is organized as follows. Section II provides an overview of related work. Section III describes the method and criteria used for performance evaluation. The experimental results are presented in Section IV. Finally, conclusions are given in Section V.

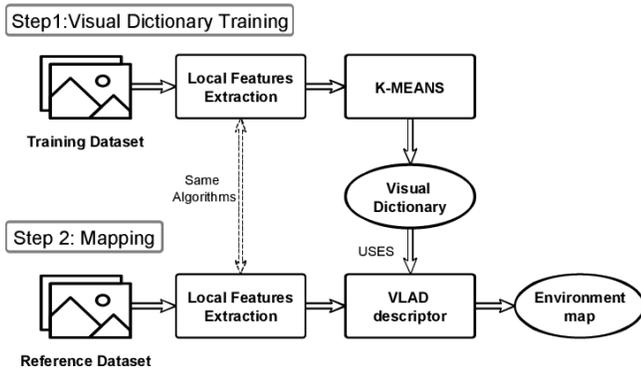


Fig. 1. Mapping occurs in two steps. First a visual dictionary to quantize the feature space is computed with k-means algorithm. In the second step, a map consisting of VLAD descriptors is built using the visual dictionary from the first step.

II. RELATED WORK

Visual place recognition is critical for developing a practical and self-reliant UAV that does not require an external tracking or positioning system [16]. The fundamental task for VPR is building and searching an image database to determine if a previously visited location is detected. The goal is to implant this ability in a UAV for it to be able to re-localize itself. However, this is challenging due to the extreme viewpoint changes experienced by a UAV and limited computational power on-board [17]. Several state-of-the-art VPR methods have been rendered unfeasible as they fail to perform satisfactorily under such conditions [2]. Convolutional Neural Networks (CNNs) display high performance for visual place recognition but are computationally intensive to be a suitable choice for a small UAV [18]. On the other hand, local feature descriptors exhibit higher potential for usage in a UAV as they are computationally less expensive. Since there is a computation-accuracy trade-off involved in the choice of a feature descriptor for UAV based visual place recognition, the selection process requires detailed analysis. To cope with the low computational power on-board a UAV, several other approaches have also been employed including the Bag of Words (BoW) approach [19]. This method can increase efficiency by creating visual dictionaries, by taking into account locally invariant feature descriptors, to match the query image vocabulary. Currently, the BoW approach tries to deal with pose change by utilizing feature descriptors, such as SIFT [11] and SURF [12] but they often fail when encountered with extreme viewpoint changes. As a result, use of range sensors [20] and structural descriptors [21] have received attention as they offer better performance to larger viewpoint changes. However, they involve higher power consumption rendering their use unfeasible for a small UAV. Other approaches undertaken include the use of FABMAP [22] as it alleviated the BoW performance but its applicability in a UAV is restricted by the high variance in performance to even small viewpoint change. Another methodology is using binary feature detectors such as ORB [15], BRISK [13] and FREAK [23] that depicted performance similar to far more

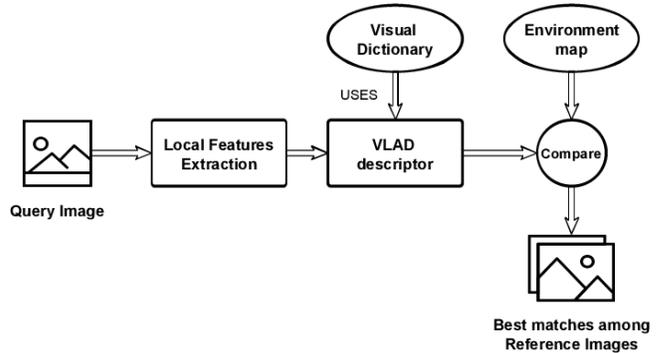


Fig. 2. The VLAD descriptor of a query image is computed and compared with those forming the environment map. Reference images are ranked by similarity between their descriptor and query image's descriptor and returned as a result.

computationally expensive features, like SIFT or SURF. With multiple possibilities, each with its perks and shortcomings, it is a tedious job to rank or select any one of these methods unbiasedly. The trade-off between VPR accuracy and the required computational effort makes the selection procedure a strenuous task. This paper proposes to explore this trade off and evaluate performance difference while avoiding any inclination towards a particular image feature descriptor by selecting a dataset comprising of both ground and aerial images.

III. EXPERIMENT SETUP

The proposed approach evaluates local image feature descriptors using VPR in the form of an image retrieval task. Local image features descriptors are used to build a map of the environment from a set of reference images showing previously visited places. During the localization phase, the map is searched to retrieve the reference images that match with a query image (i.e. a frame captured by the onboard camera). Localization succeeds if the query frame is correctly matched with the key frames in the map corresponding to the current position within the environment. The VPR algorithm and the evaluation criteria used to assess local image feature descriptors are detailed as follows.

A. Mapping

A map represents the knowledge about an environment. In our setup, a map is built from a set of images showing environment locations denoted as reference dataset, I_{REF} . For each image in I_{REF} , a set of local image features descriptors are computed and then combined into a compact image descriptor using VLAD [24], [25]. The resulting set of VLAD descriptors are finally organized in a ball-tree structure [26] in order to make the localization faster. VLAD requires a visual dictionary to quantize the feature set to combine in a single image descriptor. The visual dictionary, V_k , is computed with k-means clustering [27] using the features extracted from a training dataset, I_T . The local feature descriptor used to obtain the features set to train V_k is the same used for computing the VLAD descriptors. Figure 1 summarizes the mapping stage.

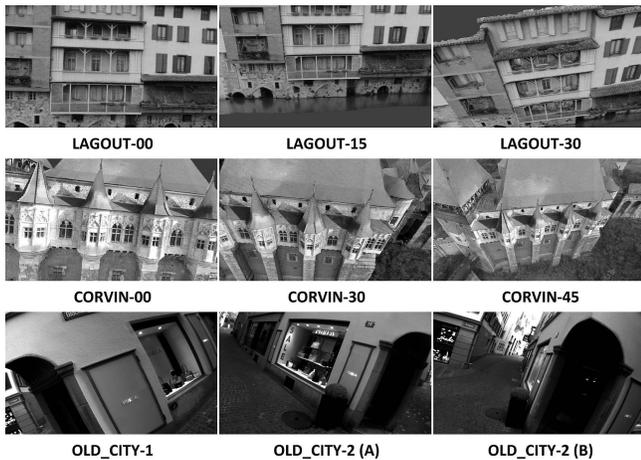


Fig. 3. A place from each dataset is shown as it appears in different loops.

B. Localization

Localization consists of finding the best matching images in I_{REF} for a query image. The image retrieval process is illustrated in the diagram in Figure 2. A VLAD descriptor is computed for the query image and compared with the VLAD descriptors of the mapped reference images. The reference images are ranked for their similarity with the input image and the resulting sequence is returned as a result of the query. Highest rank means highest similarity between a reference image and the query image. If the localization succeeds, the images ranked at top show the same place as the query image.

C. Evaluation

The highest ranked images extracted from the mapping dataset should correspond to the same place of the query image. The accuracy is evaluated with Precision and Recall (PR) while the metric for the efficiency is the required time to complete a query. In particular, the execution time includes the time spent to compute an image descriptor for the query image and the time for searching the map.

IV. RESULTS

The descriptors considered for the tests are SIFT [11], SURF [12], BRISK [13], AKAZE [14] and ORB [15]. Each descriptor has been used with its native local feature detector stage as described in the original papers. The datasets used for the tests are Lagout, Old-City and Corvin datasets [2], which include both ground and aerial images scenes captured from a wide variety of viewing angles. Lagout is a synthetic dataset consisting of 4 aerial footage captured at different angles: 0°, 15°, 30° and 45°. Corvin is similar to Lagout and includes three aerial loops around the Corvin Castle at 0°, 30° and 45°. Old-City consists of two long loops captured in an urban environment which present a wide range of views of the same locations. Figure 3 provides a sample from those datasets. The VPR algorithm used for the tests has been implemented on top of the source code available at [28] while for the local feature descriptors has been used the OpenCV [29] implementation.



Fig. 4. Scene images from VASE-JBL dataset.

The reference dataset used for mapping are the loop at 0° for Lagout and Corvin datasets and Old-City-1 for Old-City¹. The visual dictionary has been trained using a datasets which is not related with any of Lagout, Corvin and Old-City datasets. The training data consisted of an image for each of the scenes available in VASE-JBL dataset [30] for a total of 539 images showing a wide variety of outdoor and indoor scenes captured in real-world environments (Figure 4). The descriptors has been ran with the default parameter as they are set in the release 3.4.2.17 of OpenCV. The size of the visual dictionary has a significant impact on the VPR's accuracy. In order to maximize the VPR's accuracy with each descriptor, the visual dictionary length has been set at 2048 words for SURF and SIFT, 1024 for BRISK and AKAZE and 256 words for ORB.

A. Accuracy and Computation Time

Figure 5 shows the Precision-Recall curves (PR curves) for each of the assessed descriptors. The VPR exhibits the highest performance when uses SURF features with the only exception of Corvin-45 and Corvin-30 where SIFT outperform SURF by a close gap. Figure 6 shows the time required for localization for each descriptor and dataset recorded with a single thread on an Intel i7-7700K CPU. The execution is longer on Old-City because the larger number of images included in its map: there are 6711 images in Old-City-1 while only 1183 and 336 images in Corvin-00 and Lagout-00 respectively.

While SIFT could be an alternative to SURF in terms of accuracy, a VPR system based on SURF is considerably more efficient. This gap is particularly wide when the localization occurs in Old-City, where SURF allows to complete the operation in about the 60% of the time required when SIFT features are used. Although SURF based VPR can be a good trade-off in several cases, it still requires about 1s to complete the localization in Old-City using a desktop CPU. For UAVs, whose computation power is more limited, ORB can be a better option. While ORB based VPR is less accurate than SIFT and SURF based ones, it can complete the localization query 10 to 20 times faster.

B. Environment-Related Training Data

The results presented above are for a scenario where the VPR system is agnostic to the operating environment. In

¹<http://www.v4rl.ethz.ch/research/datasets-code/v4rl-wide-baseline-place-recognition-dataset.html>

Precision-Recall Curve: environment unrelated training dataset

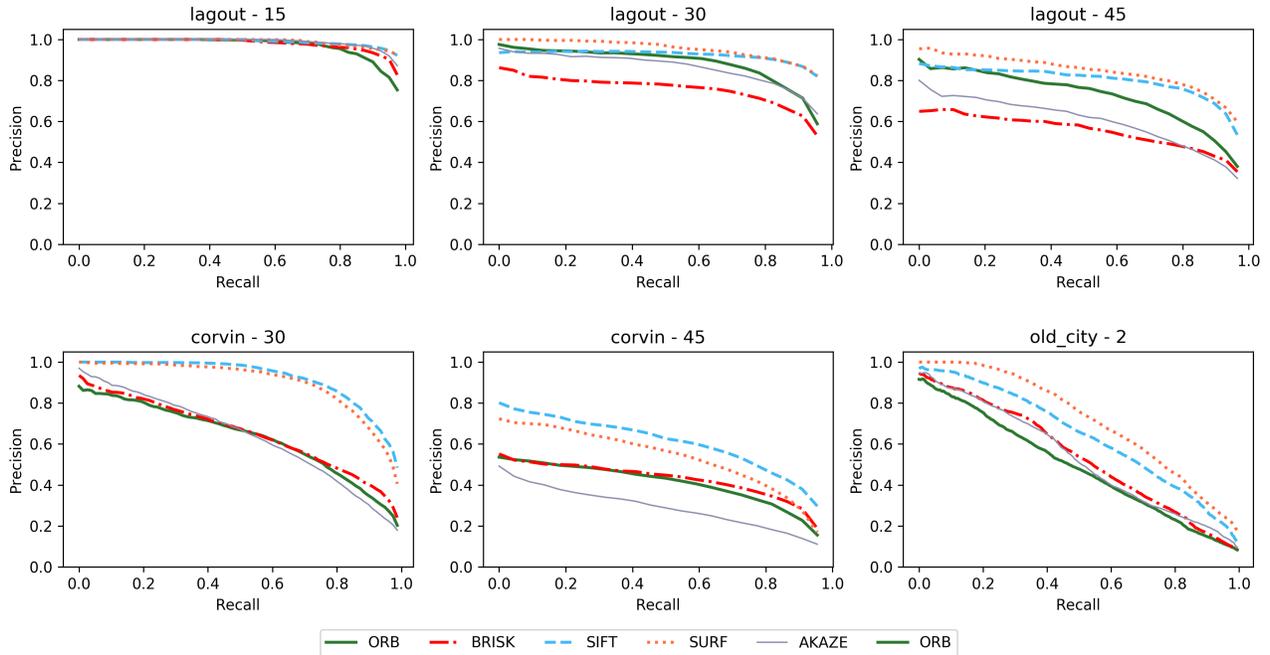


Fig. 5. Precision-Recall curves for ORB, BRISK, SIFT, SURF and AKAZE with Lagout, Corvin and Old-City datasets.

particular, the images used to train the visual dictionary is unrelated with Lagout, Corvin and Old-City datasets. This section proposes the results obtained for visual dictionaries trained with the same data used to map the environment in order to provide the VPR system with some prior knowledge of the operating environment. In details, to assess the localization performance with Lagout, the reference loop Lagout-00 has been used for both the mapping and training steps. The same setup has been used for Corvin and Old-City datasets where the reference and training data are Corvin-00 and Old-City-1 respectively.

The results confirm that the accuracy does not change significantly with respect to environment agnostic case. Figure

7 shows a comparison between the PR curves computed for the related and unrelated scenarios for Old-City dataset. None of the assessed descriptors gain a significant improvement by using environment related data to train the visual dictionary. The reason lies in the low correlation between the local features of the mapping and the test loops. Table I show the correlation coefficients from every pair of training an mapping datasets used for the experiments. The difference between the correlation coefficients of VASE-JBL and the other training dataset is neglectable with every local feature descriptor. For example, VASE-JBL on Old-City-1 are unrelated datasets and, as expected, the related correlation coefficient for AKAZE features is small (0.216). Regardless Old-City-2 and Old-City-1 show the same places, their AKAZE features can be considered unrelated as the related coefficient is just 0.214, which is very close to the value computed for VASE-JBL. Consequently, there is not a significant advantage in using images from the environment with respect to random images to train the VPR algorithm used for the tests as not any significant correlation exists between features from similar datasets.

V. CONCLUSIONS

This paper proposes a comparison of several state-of-the-art local local feature descriptors for VPR under mild to extreme viewpoint changes in small UAVs using ground-aerial image datasets. Localization accuracy is very important for VPR but it is not the only property to be considered in UAVs. As UAVs are very agile vehicles, they need to re-localize quickly but, at the same time, they are often equipped with

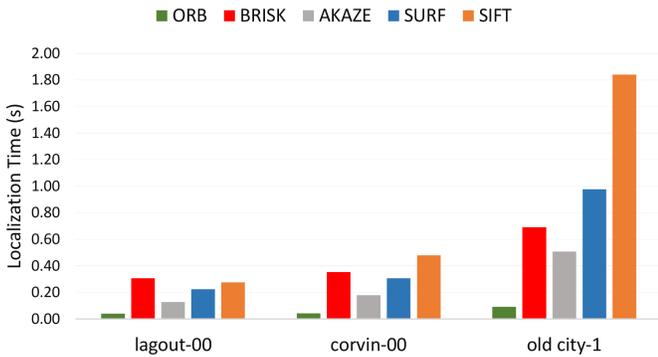


Fig. 6. The time required for localization in Lagout, Corvin and Old-City environments using an Intel i7-7700K CPU.

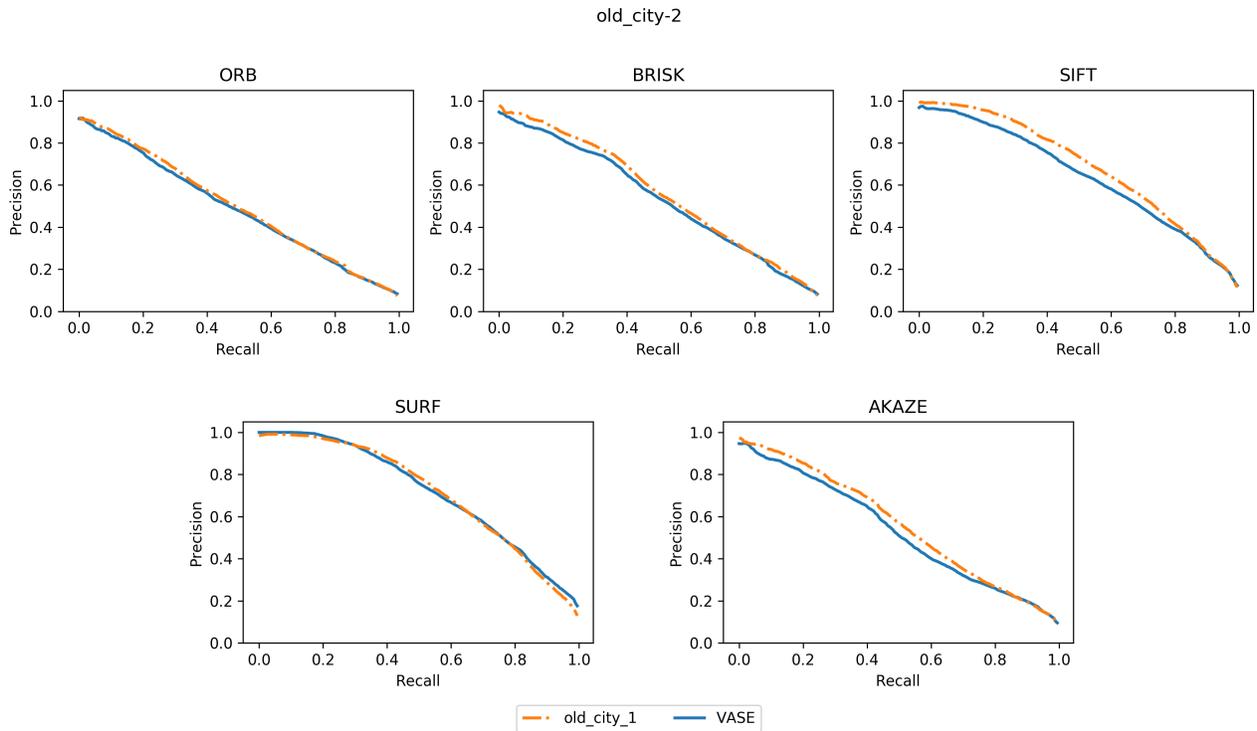


Fig. 7. A comparison between VPR trained with unrelated data from VASE-JBL dataset (blue curve) and with the same images used for the mapping (orange dot-dashed curve). The dataset used for the test is Old-City-2.

resource-constraint hardware. Driven by this consideration, the evaluation of local feature descriptors has been made on the basis of the accuracy and the computational effort to complete localization in order to determine the descriptor with the best trade-off for UAV applications. The results show that the best accuracy can be reached using SURF and SIFT descriptors at the cost of long localization time. ORB can be a better option for UAVs as it allows much faster localization while keeping a reasonable accuracy with most of the datasets considered for the experiments.

TABLE I
CORRELATION COEFFICIENTS BETWEEN TRAINING AND MAPPING DATASETS.

Training Dataset	Mapping Dataset	ORB	BRISK	SIFT	SURF	AKAZE
Lagout-15	Lagout-00	0.151	0.216	0.308	0.223	0.215
Lagout-30	Lagout-00	0.154	0.221	0.312	0.217	0.218
Lagout-45	Lagout-00	0.153	0.220	0.313	0.216	0.216
VASE-JBL	Lagout-00	0.148	0.220	0.310	0.222	0.215
Corvin-30	Corvin-00	0.153	0.216	0.313	0.220	0.210
Corvin-45	Corvin-00	0.149	0.219	0.310	0.218	0.213
VASE-JBL	Corvin-00	0.155	0.219	0.312	0.219	0.214
Old-City-1	Old-City-1	0.152	0.222	0.310	0.22	0.214
VASE-JBL	Old-City-1	0.152	0.222	0.307	0.222	0.216

REFERENCES

- [1] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Efficient deployment of multiple unmanned aerial vehicles for optimal wireless coverage," *IEEE Communications Letters*, vol. 20, no. 8, pp. 1647–1650, 2016.
- [2] F. Maffra, L. Teixeira, Z. Chen, and M. Chli, "Real-time wide-baseline place recognition using depth completion," *IEEE Robotics and Automation Letters*, 2019.
- [3] T. Villa, F. Gonzalez, B. Miljjevic, Z. Ristovski, and L. Morawska, "An overview of small unmanned aerial vehicles for air quality measurements: Present applications and future perspectives," *Sensors*, vol. 16, no. 7, p. 1072, 2016.
- [4] "Ascending Technologies," <http://www.asctec.de/en/uav-uas-drone-applications/>, accessed: 2019-03-30.
- [5] J. Li, Y. Bi, M. Lan, H. Qin, M. Shan, F. Lin, and B. M. Chen, "Real-time simultaneous localization and mapping for uav: a survey," in *Proc. of International micro air vehicle competition and conference*, 2016, pp. 237–242.
- [6] D. O. Wheeler, D. P. Koch, J. S. Jackson, T. W. McLain, and R. W. Beard, "Relative navigation: A keyframe-based approach for observable gps-degraded navigation," *IEEE Control Systems Magazine*, vol. 38, no. 4, pp. 30–48, 2018.
- [7] M. J. Milford and G. F. Wyeth, "Seqslam: Visual route-based navigation for sunny summer days and stormy winter nights," in *2012 IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 1643–1649.
- [8] N. Sünderhauf, P. Neubert, and P. Protzel, "Are we there yet? challenging seqslam on a 3000 km journey across all four seasons," in *Proc. of Workshop on Long-Term Autonomy, IEEE International Conference on Robotics and Automation (ICRA)*, 2013, p. 2013.
- [9] Z. Chen, F. Maffra, I. Sa, and M. Chli, "Only look once, mining distinctive landmarks from convnet for visual place recognition," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 9–16.
- [10] F. Maffra, Z. Chen, and M. Chli, "Viewpoint-tolerant place recognition combining 2d and 3d information for uav navigation," in *2018 IEEE*

- International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 2542–2549.
- [11] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [12] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, “Speeded-up robust features (SURF),” *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, 2008.
- [13] 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.
- [14] P. F. Alcantarilla and T. Solutions, “Fast explicit diffusion for accelerated features in nonlinear scale spaces,” *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. 34, no. 7, pp. 1281–1298, 2011.
- [15] 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.
- [16] F. Maffra, L. Teixeira, Z. Chen, and M. Chli, “Loop-closure detection in urban scenes for autonomous robot navigation,” in *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, pp. 356–364.
- [17] M. Zaffar, A. Khaliq, S. Ehsan, M. Milford, K. Alexis, and K. McDonald-Maier, “Are state-of-the-art visual place recognition techniques any good for aerial robotics?” *arXiv preprint arXiv:1904.07967*, 2019.
- [18] 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.
- [19] J. Sivic and A. Zisserman, “Video google: A text retrieval approach to object matching in videos,” in *null*. IEEE, 2003, p. 1470.
- [20] R. Dubé, D. Dugas, E. Stumm, J. Nieto, R. Siegwart, and C. Cadena, “Segmatch: Segment based place recognition in 3d point clouds,” in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 5266–5272.
- [21] T. Cieslewski, E. Stumm, A. Gawel, M. Bosse, S. Lynen, and R. Siegwart, “Point cloud descriptors for place recognition using sparse visual information,” in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 4830–4836.
- [22] 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.
- [23] A. Alahi, R. Ortiz, and P. Vanderghenst, “Freak: Fast retina keypoint,” in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. Ieee, 2012, pp. 510–517.
- [24] H. Jégou, M. Douze, C. Schmid, and P. Pérez, “Aggregating local descriptors into a compact image representation,” in *CVPR 2010-23rd IEEE Conference on Computer Vision & Pattern Recognition*. IEEE Computer Society, 2010, pp. 3304–3311.
- [25] R. Arandjelovic and A. Zisserman, “All about vlad,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2013, pp. 1578–1585.
- [26] S. M. Omohundro, *Five balltree construction algorithms*. International Computer Science Institute Berkeley, 1989.
- [27] J. A. Hartigan and M. A. Wong, “Algorithm as 136: A k-means clustering algorithm,” *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, vol. 28, no. 1, pp. 100–108, 1979.
- [28] J. Guevara Diaz, “PyVLAD,” <https://github.com/mxbi/PyVLAD>, accessed: 2019-04-04.
- [29] Itseez, “Open source computer vision library,” <https://github.com/itseez/opencv>, 2015.
- [30] S. Ehsan, A. Clark, B. Ferrarini, and K. McDonald-Maier, “Jpeg, blur and uniform light changes image database,” <http://vase.essex.ac.uk/datasets/index.html>, 2012, accessed: 2019-04-04.