

Look No Deeper: Recognizing Places from Opposing Viewpoints under Varying Scene Appearance using Single-View Depth Estimation

Sourav Garg¹, Madhu Babu V², Thanuja Dharmasiri³, Stephen Hausler¹,
Niko Suenderhauf¹, Swagat Kumar², Tom Drummond³, Michael Milford¹

Abstract—Visual place recognition (VPR) - the act of recognizing a familiar visual place - becomes difficult when there is extreme environmental appearance change or viewpoint change. Particularly challenging is the scenario where both phenomena occur simultaneously, such as when returning for the first time along a road at night that was previously traversed during the day in the opposite direction. While such problems can be solved with panoramic sensors, humans solve this problem regularly with limited field of view vision and without needing to constantly turn around. In this paper, we present a new depth- and temporal-aware visual place recognition system that solves the opposing viewpoint, extreme appearance-change visual place recognition problem. Our system performs *sequence-to-single* matching by extracting depth-filtered keypoints using a state-of-the-art depth estimation pipeline, constructing a keypoint sequence over multiple frames from the reference dataset, and comparing those keypoints to those in a single query image. We evaluate the system on a challenging benchmark dataset and show that it consistently outperforms state-of-the-art techniques. We also develop a range of diagnostic simulation experiments that characterize the contribution of depth-filtered keypoint sequences with respect to key domain parameters including degree of appearance change and camera motion.

I. INTRODUCTION

Visual Place Recognition (VPR) is a widely researched topic, with recent approaches benefiting from modern deep-learning techniques [1], [2], [3], [4], [5], especially for performance under challenging appearance variations. However, most of the existing literature dealing with such extreme appearance variations has only considered limited changes in viewpoint [6], [7]. In this paper, we address the more difficult problem of visual place recognition from opposite viewpoints under extreme appearance variations. Our system extracts depth-filtered keypoints using a state-of-the-art depth estimation pipeline [8], [9], [10], constructing a keypoint sequence over multiple frames from the reference dataset, and compares those keypoints to those in a single query image.

The color-only VPR techniques for varying environmental conditions and viewpoints have shown great promise [11], [12], [13], [4]. However, the performance achieved using these systems might be limited due to the use of only RGB information. Additional information, especially depth, is an obvious potential source for improving performance.

¹Australian Centre for Robotic Vision, Queensland University of Technology (QUT), Brisbane, Australia.

²TATA Consultancy Services, Bangalore, India.

³Australian Centre for Robotic Vision, Monash University, Melbourne, Australia.

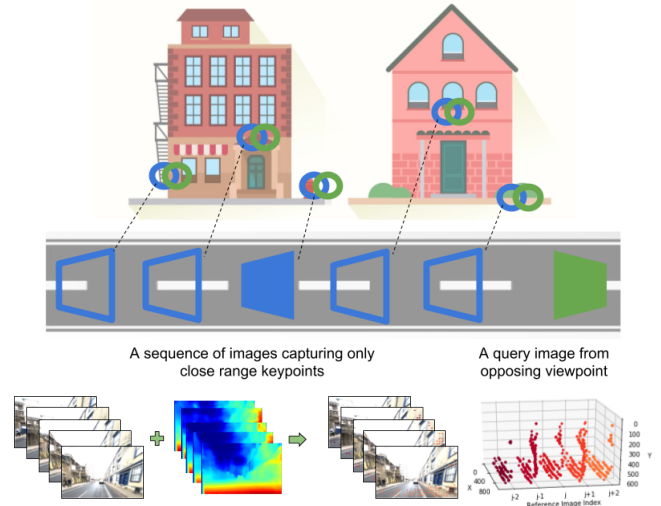


Fig. 1. *Top*: Our proposed method uses a *sequence-to-single* image matching technique to accrue depth-filtered keypoints over a sequence of reference frames (blue) to match with a single query image (green) from opposing viewpoint. *Bottom*: The reference frame sequence is combined with the depth masks to obtain a local topo-metric representation.

Moreover, depth is potentially useful for scale-aware metric pose estimation [9] though it is not the focus of this work. The use of single-view depth enables a good compromise between image space only methods [12], [11], [4] and those based on full 3D reconstruction [14], [15], where it can be challenging for the latter to deal with non-uniform appearance variations [16] and noisy depth measurements. In the same vein, our proposed system utilizes depth information within a sequence of images to create an *on-demand local topo-metric* representation [17], [18] whilst not being affected by noise in depth estimation as shown in Figure 1.

Our key contributions include:

- a novel *sequence-to-single* image matching technique¹ that exploits within-the-image information accrued over time to maximize visual overlap between images taken from opposite viewpoints,
- a depth-based keypoint filtering technique to use close-range keypoints for dealing with perceptual aliasing under extreme appearance variations,
- performance characterization with respect to depth-based keypoint filtering, reference frame sequence length, degree of appearance change, camera motion

¹The source code will be made publicly available.

and depth noise sensitivity.

The existing methods for visual place recognition that use a global image descriptor, such as those learnt using deep metric learning techniques [19], [4] encode spatial and structural features implicitly. While these methods form a strong baseline for top candidates retrieval for a query image, the explicit use of geometry can be beneficial for validating these top matching candidates [11], [20].

The problem of opposite viewpoints is often dealt with using sensor hardware solutions such as LiDARs [21] and panoramic cameras [22]. However, humans solve this problem regularly with a limited field-of-view vision and without having to constantly turn around. Under a similar constraint of limited field-of-view vision, dealing with 180 degree viewpoint change along with extreme appearance variations becomes a challenging research problem.

To address this challenge, we propose a *sequence-to-single* matching approach where a sequence of frames centered at the matching candidate location is considered for a single query image. Each of the frames in the reference candidate sequence only contributes keypoints which are within a certain depth range. As we are considering the place recognition problem for opposing viewpoints, it can be assumed that amount of visual overlap between any two matching places will only be limited to a certain depth range (on an average 30 – 40 meters for maximum visual overlap [11]). Therefore, the sequence of these depth-filtered keypoints tends to capture visual information that significantly increases the overlap with the query image. We show that the proposed technique consistently outperforms the state-of-the-art methods and also improves performance over the vanilla scenario where a single reference image is used for matching. However, the performance trends depend on the environmental conditions of the compared traverses, where extreme appearance variations tend to benefit the most due to high perceptual aliasing.

II. LITERATURE REVIEW

In visual place recognition, changing environments are a major challenge due to perceptual aliasing between the representations of places under severe appearance variations. A number of techniques have been proposed to improve the localization ability, including: leveraging temporal information [6], [23], [24], learning the appearance change over time [25], [26], [27], and extracting geometric and spatial information out of an image [28], [29], [30].

In SeqSLAM [6], utilizing sequences of recent images enables localization in challenging environments, by leveraging the scene similarity across recent images. Other sequence approaches include SMART [31], using network flows [32], and applying Conditional Random Field theory [33]. These approaches use a variety of image processing front-ends, ranging from Sum of Absolute Differences to more sophisticated methods like HOG [34] and Convolutional Neural Networks (CNNs) [35].

CNNs have demonstrated both appearance and viewpoint robustness when applied to visual place recognition [1]. In

addition to using off-the-shelf CNNs [36], the place recognition task has also been learnt end-to-end [4]. In HybridNet, an off-the-shelf CNN is improved by training on appearance-change images [4], [27]. In HybridNet [27], pyramid pooling is used to improve the viewpoint robustness, however this pooling method is still defined by specific regions within a feature map. As an improvement, deep-learnt features have been intelligently extracted out of keypoint locations within these feature maps [37]. [38] proposes cross-convolutional pooling, by pooling features using the spatial position of activations in the subsequent layer. In [3], keypoint locations are extracted out of a late convolutional layer, whereas LoST [11] uses semantic information to represent places and obtain keypoint locations. LoST improves the localization ability across viewpoint variations as extreme as front to rear-view imagery. However, it relies mostly on visual semantics and spatial layout verification to achieve high performance; in this work, we delve deeper into the efficacy of CNN-based keypoints and descriptors for VPR under challenging scenarios by using depth- and temporal-aware system.

The use of geometric information for VPR has been shown to improve performance in recent works [12], [11]. In the same vein, we explore the use of depth information for improving VPR performance under vast variations in appearance and viewpoint. Depth estimation is best achieved using stereoscopic images [39], however recent advances have enabled depth estimation from monocular images. Geometric constraints [40] and non-parametric sampling [41] have been used to extract depth out of a monocular image. Improved performance can be gained by training a CNN to estimate pixel-by-pixel depth within an image [42]. In recent years, several deep networks have been proposed to estimate per-pixel-depth map and visual odometry. For instance, [43] is one of the earliest works in predicting depth by regressing with the ground truth for monocular images. The authors in [30] additionally use the motion parallax between image pairs, enabling robust depth estimation for novel scenes. Rather than using supervised learning to estimate depth, [44], [8] use unsupervised learning in an end-to-end framework. CNN-based depth estimation has been employed for dense monocular slam and semantic 3D reconstruction [45]. Our proposed approach exploits depth estimates to filter keypoints from a sequence of reference frames that are beyond a certain range for the expected visual overlap between front- and rear-view images. This forms a local topo-metric representation of the reference candidate match location which is a good compromise between full 3D reconstruction and single image only methods. However, methods like [45] will be complementary to our proposed approach.

III. PROPOSED APPROACH

We use a hierarchical place matching pipeline [11] where top-N candidate matches are generated using cosine distance based descriptor comparison of a query image with the reference database using state-of-the-art retrieval methods [4], [11], [46], [38]. The query image is then matched with these top-N candidates using depth-filtered keypoints extracted

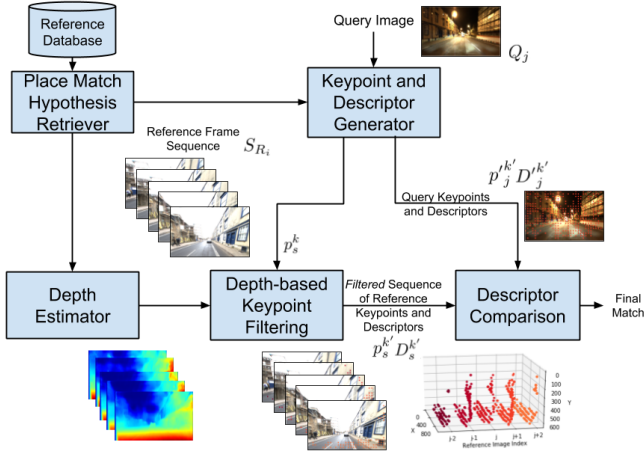


Fig. 2. Flowchart of our proposed approach.

from a sequence of reference frames centered at the candidate index. Figure 2 shows a flow diagram of our proposed approach as explained in subsequent sections.

A. Keypoint and Descriptor Extraction

For a given pair of matching images, we extract keypoints and descriptors from the *conv5* layer of the ResNet101 [47] as described in details in [11] and briefly here for sake of completion. For a *conv5* tensor of dimension $W \times H \times C$, representing width, height, and the number of channels (or feature maps) of the tensor, a keypoint location is determined for each of the channels based on the maximally-activated 2D location within that channel. Using the same tensor, a descriptor at the keypoint location is extracted as a C dimensional vector along the third axis. As only one keypoint is extracted per channel, we obtain a total of C such keypoints and descriptors. For a given pair of images, the correspondences between the keypoints are assumed based on their channel index.

B. Unsupervised Monocular Depth Estimation

Given the reference database images, a per-pixel-depth map is estimated using the framework described in UnDE-MoN [8]. The network estimates disparity maps which are then converted into depth maps using the camera parameters of the training dataset. For a keypoint k at a pixel location p_i^k , depth $u_{p_i^k}$ is estimated from the disparity $u'_{p_i^k}$ using the formula $u_{p_i^k} = \frac{fb}{u'_{p_i^k}}$ where f is the focal length and b is the baseline distance between the stereo image pairs from the KITTI dataset [48] used to train the model. In this work, we directly used the KITTI trained models from [8] without any fine tuning performed on the datasets used in our experiments.

C. Depth-Filtered Keypoints Sequence

For a query image, Q_j and a top matching reference candidate R_i , keypoints $p_j^{k'}$ and p_i^k and descriptors $D_j^{k'}$ and D_i^k are extracted respectively. Further, a sequence of reference frames S_{R_i} of length l is considered across R_i ,

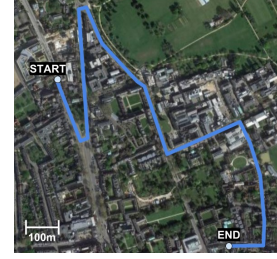


Fig. 3. Aerial view of ground truth trajectories for Oxford Robotcar dataset. Source: Google Map

that is, $S_{R_i} = [R_{i-l/2}, R_{i+l/2}]$ as shown in Figure 2. Then, a depth range threshold d is used to filter out the keypoints that are far from the camera. Therefore, a set of keypoints is obtained from the reference sequence as below:

$$p_s^{k'} = \{p_s^k \mid u_{p_s^k} < d\} \quad \forall k \in C, s \in S_{R_i} \quad (1)$$

where k' is a filtered keypoint index within $C' \subseteq C$.

D. Descriptor comparison

For each of these keypoints $p_s^{k'}$ in the reference sequence, there exists a corresponding keypoint in the query image². We use cosine distance between the descriptors of these corresponding keypoints to obtain minimum descriptor distance from the query keypoint to the corresponding keypoint in any of the reference sequence frames:

$$r_{ji}^{k'} = \min_s \left(1 - \frac{D_j^{k'} \cdot D_s^{k'}}{\|D_j^{k'}\|_2 \|D_s^{k'}\|_2} \right) \quad \forall k' \in C', s \in S_{R_i} \quad (2)$$

An average over all the filtered keypoints' descriptor distances $r_{ji}^{k'}$ gives the matching score r_{ji} between the query image and the reference candidate.

The least scoring candidate is then considered as the final image match for the query.

IV. EXPERIMENTAL SETUP

A. Oxford Robotcar Dataset

The Oxford Robotcar Dataset [49] comprises traverses of Oxford city during different seasons, time of day and weather conditions, capturing images using cameras pointing in all four directions. We used an initial 2.5 km traverse from front- and rear-view cameras for four different environmental settings: Overcast Autumn, Night Autumn, Dawn Winter and Overcast Summer³. We used the GPS data to sample image frames at a constant distance of approximately 2 meters.

²The correspondences are based on the channel (feature map) index (k') within the *conv5* tensor

³2014-12-09-13-21-02, 2014-12-10-18-10-50, 2015-02-03-08-45-10 and 2015-05-19-14-06-38 respectively in [49]

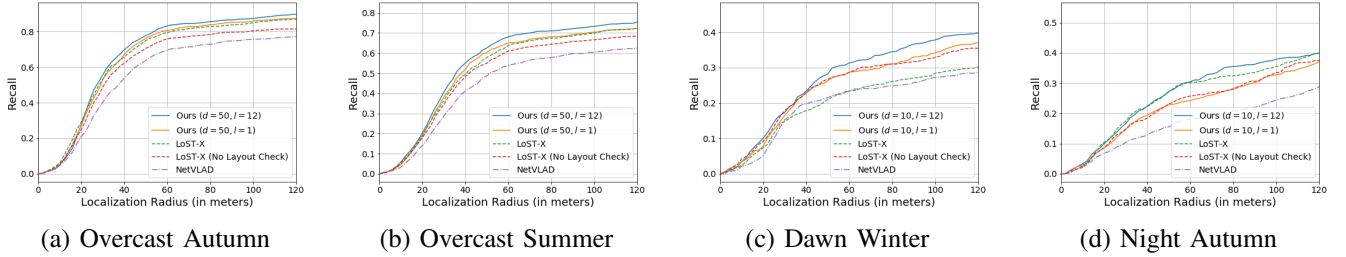


Fig. 4. Performance comparison against state-of-the-art methods for opposite viewpoints and changing conditions. These front-view imagery from the four traverses: (a) Overcast Autumn, (b) Overcast Summer, (c) Dawn Winter, and (d) Night Autumn was compared against the rear-view imagery from Overcast Autumn traverse.

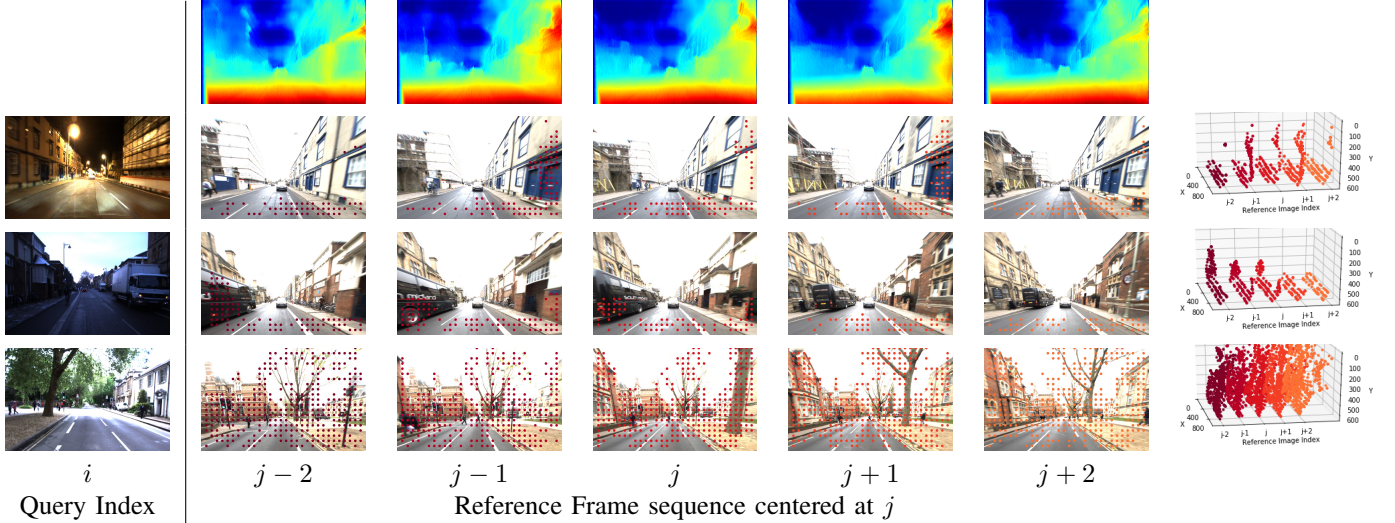


Fig. 5. Example Matches: Keypoints from the front-view query images (*leftmost column*) are matched against a reference frame sequence using only depth-filtered keypoints (marked in different colors) accrued over the sequence. The *rightmost column* shows the sequential collection of these keypoints such that the horizontal axis topologically connects the sequence of images, however, for each image metric depth is used to separate the points. The *topmost row* shows depth masks for the reference image sequence used in the second row.

B. Ground Truth and Evaluation Method

We use GPS information to generate place recognition ground truth for the Oxford traverses. However, unlike the forward-forward image matching scenario, the opposite viewpoint place matching involves a *visual offset* of approximately 30–40 meters where the visual overlap between the matching pair of images becomes maximum, as described in [11]. We use recall as a performance measure, that is equal to the number of true positives detected divided by the total number of positives in the dataset comparison. A match is considered a true positive if it lies within a certain radius of its ground truth GPS location and this radius is varied to generate the performance curves.

C. Comparative Study

We compare our proposed approach with two state-of-the-art methods: ‘NetVLAD’ [4] and ‘LoST-X’ [11]. LoST-X uses spatial layout matching of semantically-filtered keypoints; we also include the ‘LoST-X (No Layout Check)’ version in order to compare against semantic filtering based *conv5* descriptor matching. For our proposed approach, we show results for two scenarios: one with the reference frame sequence length l set to 1 which means traditional single-to-single image matching, and the other with l set to 12

which is approximately 25 meters for the Oxford dataset. The depth range threshold is set to 50 and 10 meters for similar and different environmental conditions respectively. The choice of parameters is primarily based on the previous study [11] that shows 30–40 meters to be an average distance between two images from opposing viewpoints leading to the maximum visual overlap. Further, in Section V-B, we show performance characteristics of our proposed system with respect to these parameters that support choosing depth threshold depending on the degree of appearance variations. For a fair comparison, for all the methods, we use top 5 match hypothesis generated by cosine distance comparison of NetVLAD descriptors.

V. RESULTS

A. Performance comparison

Figure 4 shows performance comparison of our proposed approach against state-of-the-art methods for opposite viewpoint image matching under varying appearance of the environment. It can be observed that our proposed method consistently outperforms both NetVLAD and LoST-X. The performance gains are even more when compared against LoST-X without employing its spatial layout verification.

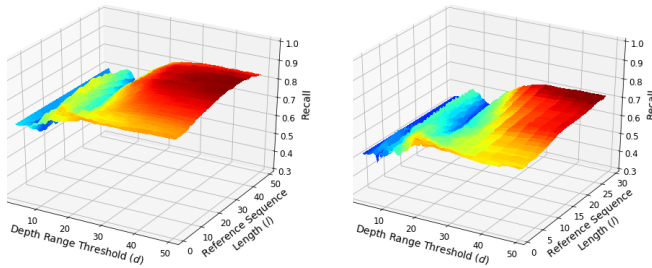


Fig. 6. *Similar environmental conditions*: The performance curves comparing Autumn Day Rear with Autumn Day Front (*left*) and (b) Summer Day Front (*right*) show that performance tends to increase as more visual information becomes available by considering distant points from the camera using depth information and a longer sequence of reference frames.

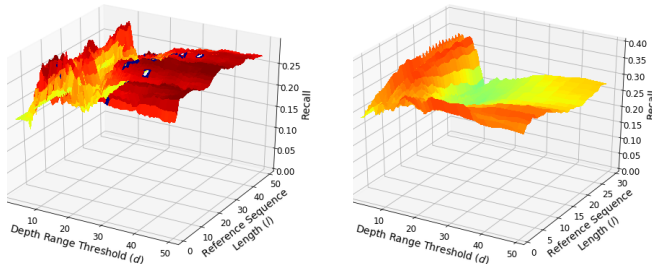


Fig. 7. *Different environmental conditions*: The performance curves comparing Autumn Day Rear with Night Front (*left*) and Winter Front (*right*) show that distant points from the camera tend to decrease the performance. Further, longer reference frame sequence does not provide any significant performance gains unlike those observed when environmental conditions remained similar as shown in Figure 6.

This indicates that spatial layout matching is also a key step to attain higher performance and can potentially complement our proposed approach to improve performance even further. The performance curves also show consistently that using only a single reference image based matching has sub-optimal performance than the reference frame sequence based matching.

Qualitative Results: Figure 5 shows qualitative matches for opposite viewpoint image matching under varying environmental conditions. A sequence of frames centered at the candidate match location (index= j) are shown with their keypoints within a certain depth range. The *rightmost column* shows the sequential collection of these keypoints such that the horizontal axis topologically connects the sequence of images, however, for each image metric depth is used to separate the points. The *topmost row* shows depth masks for the reference image sequence used in the second row. It can be observed that due to imperfections in the estimated depth, filtering of keypoints within a given depth range is not always consistent, however, due to visual overlap between consecutive frames, this doesn't pose a problem.

B. Performance Characteristics

In this section, we show the performance characteristics of our proposed system with respect to the two system parameters: depth range threshold, d and reference frame sequence length, l . These characterizations are discussed below with respect to the extent of appearance variations

due to varying environmental conditions. Further, we show system's characteristics with respect to different camera speed and the quality of depth estimation.

1) *Moderate Appearance Variations*: Figure 6 shows that for opposite viewpoints under moderate or no changes in environmental conditions, more visual information in the form of longer reference sequence length and larger pool of keypoints within each of the reference frames, helps attain high performance. This shows that descriptors used from within the CNN are robust to moderate variations in appearance and are able to discriminate between a large number of false and true keypoint correspondences, therefore, able to utilize additional visual information.

2) *Extreme Appearance Variations*: Figure 7 shows the performance characteristics for opposite viewpoints under extreme appearance variations from day to night and autumn to winter. It can be observed that, unlike the previous scenario of moderate appearance variations, performance tends to decrease as the visual information tends to increase, especially by allowing more keypoints per reference frame. This can be attributed to high perceptual aliasing caused due to extreme appearance variations. Therefore, accumulating keypoints within a short range of camera over the reference frame sequence helps attain optimal performance. Further within the optimal performance region, the reference frame sequence length tends to have limited effect on the performance, therefore, allowing the selection of a shorter sequence length.

3) *Camera Speed*: Figure 8 shows variation of performance with respect to changing camera speed for the Day-Front and Night-Rear comparison. This is simulated by skipping frames within the reference frame sequence at: (a) $1\times$, (b) $2\times$, and (c) $4\times$. It can be observed that with a higher camera speed (or frame skip rate), the peak performance decreases, however, the effect of using a shorter depth range becomes more prominent.

4) *Quality of Depth Estimation*: We use Synthia dataset [50] for evaluating the effect of quality of depth⁴ for our proposed approach. For this purpose, we use the front- and rear-view images from the Dawn and Fall traverses of Sequence-02 that have very different environmental conditions, as also utilized in [11]. Figure 9 shows performance comparison using (a) ground truth depth and (b) estimated depth. It can be observed that the performance trends remain similar to those observed for Oxford dataset for opposite viewpoints and varying appearance. Further, with the use of ground truth depth, peak performance is slightly higher and performance variations are smooth with respect to the system parameters. The bottom row in Figure 9 also shows a comparison between the ground truth and estimated depth mask for an input image from Synthia reference traverse.

VI. DISCUSSION

The performance characterizations in Figure 6 and 7 show that the use of more visual information by including more

⁴The depth estimation network used in this paper [8] is trained on KITTI dataset [48].

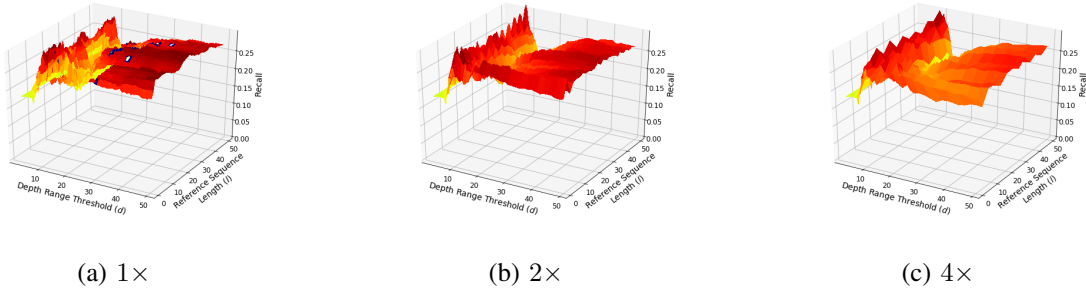


Fig. 8. *Effects of camera speed*: The performance trends tend to remain similar even when the reference frame sequence is considered at a different camera speed (simulated by skipping frames in the sequence): (a) $1\times$, (b) $2\times$, and (c) $4\times$. It can also be observed that with a higher camera speed (frame skip rate), the peak performance decreases, however, the effect of using a shorter depth range becomes more prominent.

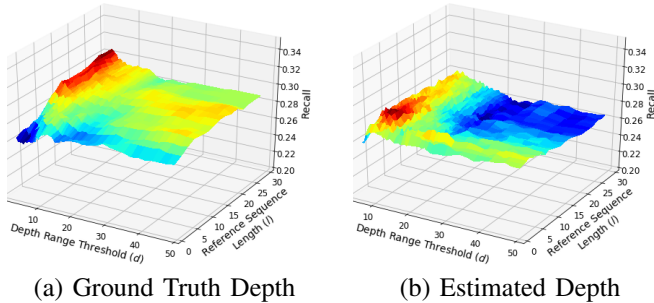


Fig. 9. *Quality of Depth*: *Top*: Performance trends tend to remain similar for both (a) Ground truth depth and (b) Estimated depth. However, peak performance is slightly better when using ground truth depth along with a smooth variation with respect to the system parameters. *Bottom*: Ground truth and estimated depth masks for an input image.

Input Image Ground truth Depth Estimated Depth

distant points and using longer reference frame sequence has varied implications, depending on the variations in the scene appearance. While similar environmental settings tend to benefit from increasing visual information, performance drops drastically when images are matched across different environmental conditions. This drop in performance occurs because of perceptual aliasing among large number of keypoints matched under extreme appearance variations. Our proposed approach allows control of this visual information to achieve optimal performance by using close-range keypoints collected over multiple reference frames.

Most of the existing literature for VPR that deals with challenging appearance variations, generally deals with similar or moderate variations in viewpoints. In such cases, the above finding often gets discounted because the overall structure of the scene remains mostly similar that massively aids in matching, especially when images are down-sampled, for example, patch-normalization in SeqSLAM [6], HoG in [7], and flattened deep-learned tensors in [51], [52], [46], [36]. Furthermore, as we use the deep metric learning method NetVLAD [4] as a baseline to generate top candidate

matches, it shows that such methods, despite implicitly encoding the spatial and structural information, still have a large room for improvement. The above analysis helps in gaining insights about the behavior of deep-learned convolutional descriptors under different appearance conditions when no assumptions are made about the overall structural similarity.

VII. CONCLUSION AND FUTURE WORK

Visual place recognition under opposing viewpoints and varying environmental conditions is a challenging problem and requires effective use of both scene appearance and scene geometry. In this paper, we proposed a *sequence-to-single* matching approach where an on-demand local topometric representation of reference image sequence was used to match with a single query image using depth-based keypoint filtering. We showed that our proposed system performs better than the state of the art on challenging benchmark datasets. Further, the performance characterizations also revealed that the amount of visual information can be controlled using depth-based filtering to reduce perceptual aliasing, thereby leading to optimal performance under extreme appearance variations. In future, we plan to extend our work to a visual SLAM pipeline for robust metric localization and mapping under extreme appearance and viewpoint variations.

REFERENCES

- [1] N. Sunderhauf, S. Shirazi, A. Jacobson, F. Dayoub, E. Pepperell, B. Upercroft, and M. Milford, "Place recognition with convnet landmarks: Viewpoint-robust, condition-robust, training-free," *Proceedings of Robotics: Science and Systems XII*, 2015.
- [2] O. Vysotska and C. Stachniss, "Lazy data association for image sequences matching under substantial appearance changes," *IEEE Robotics and Automation Letters*, vol. 1, no. 1, pp. 213–220, 2016.
- [3] Z. Chen, F. Maffra, I. Sa, and M. Chli, "Only look once, mining distinctive landmarks from convnet for visual place recognition," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 9–16.
- [4] 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.
- [5] H. Noh, A. Araujo, J. Sim, T. Weyand, and B. Han, "Large-scale image retrieval with attentive deep local features," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3456–3465.

- [6] M. J. Milford and G. F. Wyeth, "Seqslam: Visual route-based navigation for sunny summer days and stormy winter nights," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 1643–1649.
- [7] T. Naseer, L. Spinello, W. Burgard, and C. Stachniss, "Robust visual robot localization across seasons using network flows," in *Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014.
- [8] V. Babu, A. Majumder, K. Das, S. Kumar, et al., "A deeper insight into the undemon: Unsupervised deep network for depth and ego-motion estimation," *arXiv preprint arXiv:1809.00969*, 2018.
- [9] T. Dharmasiri, A. Spek, and T. Drummond, "Eng: End-to-end neural geometry for robust depth and pose estimation using cnns," *arXiv preprint arXiv:1807.05705*, 2018.
- [10] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox, "Demon: Depth and motion network for learning monocular stereo," in *IEEE Conference on computer vision and pattern recognition (CVPR)*, vol. 5, 2017, p. 6.
- [11] S. Garg, N. Suenderhauf, and M. Milford, "Lost? appearance-invariant place recognition for opposite viewpoints using visual semantics," 2018.
- [12] A. Gawel, C. Del Don, R. Siegwart, J. Nieto, and C. Cadena, "X-view: Graph-based semantic multi-view localization," *arXiv preprint arXiv:1709.09905*, 2017.
- [13] T. Naseer, G. L. Oliveira, T. Brox, and W. Burgard, "Semantics-aware visual localization under challenging perceptual conditions," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2017.
- [14] 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.
- [15] J. Engel, T. Schöps, and D. Cremers, "Lsd-slam: Large-scale direct monocular slam," in *European Conference on Computer Vision*. Springer, 2014, pp. 834–849.
- [16] H. Alismail, B. Browning, and S. Lucey, "Direct visual odometry using bit-planes," *arXiv preprint arXiv:1604.00990*, 2016.
- [17] H. Badino, D. Huber, and T. Kanade, "Visual topometric localization," in *Intelligent Vehicles Symposium (IV), 2011 IEEE*. IEEE, 2011, pp. 794–799.
- [18] S. Bazeille and D. Filliat, "Incremental topo-metric slam using vision and robot odometry," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 4067–4073.
- [19] M. A. U. G. H. Lee, "Pointnetvlad: Deep point cloud based retrieval for large-scale place recognition," in *IEEE Conference on computer vision and pattern recognition (CVPR)*, 2018.
- [20] S. Garg, N. Suenderhauf, and M. Milford, "Don't look back: Robustifying place categorization for viewpoint- and condition-invariant place recognition," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- [21] R. W. Wolcott and R. M. Eustice, "Robust lidar localization using multiresolution gaussian mixture maps for autonomous driving," *The International Journal of Robotics Research*, vol. 36, no. 3, pp. 292–319, 2017.
- [22] R. Arroyo, P. F. Alcantarilla, L. M. Bergasa, J. J. Yebes, and S. Gámez, "Bidirectional loop closure detection on panoramas for visual navigation," in *Intelligent Vehicles Symposium Proceedings, 2014 IEEE*. IEEE, 2014, pp. 1378–1383.
- [23] Y. Liu and H. Zhang, "Towards improving the efficiency of sequence-based slam," in *Mechatronics and Automation (ICMA), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1261–1266.
- [24] P. Hansen and B. Browning, "Visual place recognition using hmm sequence matching," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*. IEEE, 2014, pp. 4549–4555.
- [25] P. Neubert, N. Sunderhauf, and P. Protzel, "Appearance change prediction for long-term navigation across seasons," in *Mobile Robots (ECMR), 2013 European Conference on*. IEEE, 2013, pp. 198–203.
- [26] S. M. Lowry, M. J. Milford, and G. F. Wyeth, "Transforming morning to afternoon using linear regression techniques," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 3950–3955.
- [27] Z. Chen, A. Jacobson, N. Sunderhauf, B. Upcroft, L. Liu, C. Shen, I. Reid, and M. Milford, "Deep learning features at scale for visual place recognition," *arXiv preprint arXiv:1701.05105*, 2017.
- [28] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*, 2007, pp. 225–234.
- [29] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [30] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox, "Demon: Depth and motion network for learning monocular stereo," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5622–5631.
- [31] E. Pepperell, P. I. Corke, and M. J. Milford, "All-environment visual place recognition with smart," in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 1612–1618.
- [32] T. Naseer, W. Burgard, and C. Stachniss, "Robust visual localization across seasons," *IEEE Transactions on Robotics*, 2018.
- [33] C. Cadena, D. Gálvez-López, J. D. Tardós, and J. Neira, "Robust place recognition with stereo sequences," *IEEE Transactions on Robotics*, vol. 28, no. 4, pp. 871–885, 2012.
- [34] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 1. IEEE, 2005, pp. 886–893.
- [35] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [36] N. Sünderhauf, S. Shirazi, F. Dayoub, B. Upcroft, and M. Milford, "On the performance of convnet features for place recognition," in *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE, 2015, pp. 4297–4304.
- [37] N. Ufer and B. Ommer, "Deep semantic feature matching," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 5929–5938.
- [38] L. Liu, C. Shen, and A. van den Hengel, "The treasure beneath convolutional layers: Cross-convolutional-layer pooling for image classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4749–4757.
- [39] P. R. Induchoodan, M. J. Josemartin, and P. R. Geetharanjin, "Depth recovery from stereo images," in *2014 International Conference on Contemporary Computing and Informatics (IC3I)*, 2014, pp. 745–750.
- [40] M. Reza, W. Martin, and A. Anelia, "Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints," *arXiv preprint arXiv:1706.07144*, 2018.
- [41] K. Karsch, C. Liu, and S. B. Kang, "Depth transfer: Depth extraction from video using non-parametric sampling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 11, pp. 2144–2158, 2014.
- [42] F. Liu, C. Shen, G. Lin, and I. Reid, "Learning depth from single monocular images using deep convolutional neural fields," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 10, pp. 2024–2039, 2016.
- [43] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, ser. NIPS'14. Cambridge, MA, USA: MIT Press, 2014, pp. 2366–2374. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2969033.2969091>
- [44] R. Garg, V. K. BG, G. Carneiro, and I. Reid, "Unsupervised cnn for single view depth estimation: Geometry to the rescue," in *European Conference on Computer Vision*. Springer, 2016, pp. 740–756.
- [45] K. Tateno, F. Tombari, I. Laina, and N. Navab, "CNN-SLAM: real-time dense monocular SLAM with learned depth prediction," *CoRR*, vol. abs/1704.03489, 2017. [Online]. Available: <http://arxiv.org/abs/1704.03489>
- [46] A. Babenko and V. Lempitsky, "Aggregating local deep features for image retrieval," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1269–1277.
- [47] G. Lin, A. Milan, C. Shen, and I. Reid, "Refinenet: Multi-path refinement networks with identity mappings for high-resolution semantic segmentation," *arXiv preprint arXiv:1611.06612*, 2016.
- [48] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, June 2012, pp. 3354–3361.
- [49] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 year, 1000 km: The oxford robotcar dataset," *IJ Robotics Res.*, vol. 36, no. 1, pp. 3–15, 2017.

- [50] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3234–3243.
- [51] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "Cnn features off-the-shelf: an astounding baseline for recognition," in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2014 IEEE Conference on*. IEEE, 2014, pp. 512–519.
- [52] Z. Chen, O. Lam, A. Jacobson, and M. Milford, "Convolutional neural network-based place recognition," *arXiv preprint arXiv:1411.1509*, 2014.