

2D3D-MatchNet: Learning to Match Keypoints Across 2D Image and 3D Point Cloud

Mengdan Feng¹, Sixing Hu², Marcelo H Ang Jr¹ and Gim Hee Lee²

Abstract—Large-scale point cloud generated from 3D sensors is more accurate than its image-based counterpart. However, it is seldom used in visual pose estimation due to the difficulty in obtaining 2D-3D image to point cloud correspondences. In this paper, we propose the 2D3D-MatchNet - an end-to-end deep network architecture to jointly learn the descriptors for 2D and 3D keypoint from image and point cloud, respectively. As a result, we are able to directly match and establish 2D-3D correspondences from the query image and 3D point cloud reference map for visual pose estimation. We create our *Oxford 2D-3D Patches dataset* from the Oxford Robotcar dataset with the ground truth camera poses and 2D-3D image to point cloud correspondences for training and testing the deep network. Experimental results verify the feasibility of our approach.

I. INTRODUCTION

Visual pose estimation refers to the problem of estimating the camera pose with respect to the coordinate frame of a given reference 3D point cloud map. It is the foundation of visual Simultaneous Localization and Mapping (vSLAM) [1], [2], and Structure-from-Motion (SfM) [3], which are extremely important to applications such as autonomous driving [4] and augmented reality [5] etc. A two-step approach [3] is commonly used for visual pose estimation - (1) establish 2D-3D keypoint correspondences between the 2D image and 3D reference map, and (2) apply the Perspective-n-Point (PnP) [6] algorithm to compute the camera pose with respect to the coordinate frame of the 3D reference map with at least three 2D-3D correspondences. The 3D point cloud of the reference map is usually built from a collection of images using SfM, and the associated keypoint descriptors, e.g. SIFT [7], are stored with the map to facilitate the establishment of 2D-3D correspondences in visual pose estimation.

It is eminent that the accuracy of visual pose estimation strongly depends on the quality of the 3D reference map. Unfortunately, it is hard to ensure the quality of the 3D point cloud reconstructed from SfM, since most images are taken by close-to-market photo-sensors that are noisy. Furthermore, absolute scale of the reconstructed 3D point cloud is not directly available from SfM and has to be obtained from other sources. The need to store keypoint descriptors from multiple views of the same keypoint also increases the memory consumption of the map. In contrast, a 3D point cloud map built from Lidars has the advantages of higher accuracy and absolute scale is directly observed.

Despite the advantages, Lidars are seldom used to build the 3D reference map because of the lack of a descriptor that allows direct matching of the keypoints extracted from a 2D image and 3D point cloud.

In this paper, we propose the 2D3D-MatchNet - a novel deep network approach to jointly learn the keypoint descriptors of the 2D and 3D keypoints extracted from an image and a point cloud. We use the existing detectors from SIFT [7] and ISS [8] to extract the keypoints of the image and point cloud, respectively. Similar to most deep learning methods, an image patch is used to represent an image keypoint, and a local point cloud volume is used to represent a 3D keypoint. We propose a triplet-like deep network to concurrently learn the keypoint descriptors of a given image patch and point cloud volume such that the distance in the descriptor space is small if the 2D and 3D keypoint are a matching pair, and large otherwise. The descriptors of the keypoints from both the image and point cloud are generated through our trained network during inference. The EPnP [9] algorithm is used to compute the camera pose based on the 2D-3D correspondences. We create our *Oxford 2D-3D Patches dataset* with 432,982 pairs of matching 2D-3D keypoints based on the Oxford dataset. We conduct extensive experiments on our dataset to verify that sufficient inliers for visual pose estimation can be found based on our learned feature vectors. Pose estimation succeeds without any prior on 2D-3D correspondences.

Contributions (1) To the best of our knowledge, we are the first to propose a deep learning approach to learn the descriptors that allow direct matching of keypoints across a 2D image and 3D point cloud. (2) Our approach makes it possible for the use of Lidars to build more accurate 3D reference map for visual pose estimation. (3) We create a dataset with huge collection of 2D-3D image patch to 3D point cloud volume correspondences, which can be used to train and validate the network.

II. RELATED WORK

Traditional localization Most existing work on visual pose estimation can be classified into two categories: (1) local structure based methods and (2) global appearance based methods. 2D-3D correspondences are first established from SIFT features given the query image and 3D scene model [10], [11] for local structure based methods. Each local feature pair votes for its own pose independently, without considering other pairs in the image. Then a minimal solver algorithm, e.g. [9], combined with RANSAC iterations, is used for robust pose estimation. Global appearance based

¹Department of Mechanical Engineering, National University of Singapore fengmengdan@u.nus.edu, mpeangh@nus.edu.sg

²Computer Vision and Robotic Perception (CVRP) Lab, Department of Computer Science, National University of Singapore hu.sixing@u.nus.edu, gimhee.lee@nus.edu.sg

localization methods [12], [13] aggregate all local features of the query image to a global descriptor and localize the query image by matching against its nearest neighbour in current image database as a retrieval problem.

Learnable localization Deep learning is increasingly applied to the visual pose estimation since learned features are shown to be more robust against environmental changes, e.g. lighting and weather changes, compared with methods based on hand-crafted features such as SIFT [14]. Existing deep network models solve the localization problem from different aspects. [14], [15] learn to regress the 6D camera pose directly based on single image. [16], [17] learn to predict pixel-wise image to scene correspondences, followed by a RANSAC-like optimization algorithm for pose estimation. [18] proposes to localize an image by learnable probabilistic 2D-3D correspondences and iterative refinement of camera pose. However, they cannot generalize to unseen environment due to the global pose estimation mechanism. [19] proposes to learn robust 3D point cloud descriptors by fusing semantic and geometric information for long-term localization.

Deep similarity learning Deep similarity learning is widely used to achieve the information retrieval task. Two common architectures of deep similarity learning are the Siamese network and the triplet network. The Siamese network learns the similarity relationship between a pair of inputs [20], [21]. Most existing work shows better results on the Triplet architecture [22]–[26]. Liao *et. al* [23] use the triplet network to re-identify a person’s identity. Vo and Hays [24] conduct experiments on the ground-to-aerial image retrieval task using both Siamese architecture and Triplet architecture. They show that Triplet architecture performs better. The Triplet network outperforms the Siamese network because it can jointly pull the positive sample to the anchor while pushing the negative sample away. In our work, our proposed network is based on the Triplet network.

III. APPROACH

In this section, we outline our pipeline for visual pose estimation with a 2D query image and 3D point cloud reference map built from Lidar scans. We first introduce the overview of our pipeline in section III-A. In section III-B, we describe our novel 2D3D-MatchNet - a deep network to jointly extract the descriptors of the 2D and 3D keypoints from an image and a point cloud. The training loss is given in section III-B. Finally, we discuss the pose estimation algorithm we use to compute the camera pose given at least three 2D-3D correspondences in section III-C.

A. Overview

Given a query image I and the 3D point cloud reference map M of the scene, the objective of visual pose estimation is to compute the absolute camera pose $P = [R | t]$ of the query image I with respect to the coordinate frame of the 3D point cloud reference map M . Unlike existing visual pose methods which associate image-based descriptors, e.g. SIFT [7], to each 3D point in the reference map, we propose the 2D3D-MatchNet - a deep network to jointly learn the descriptors

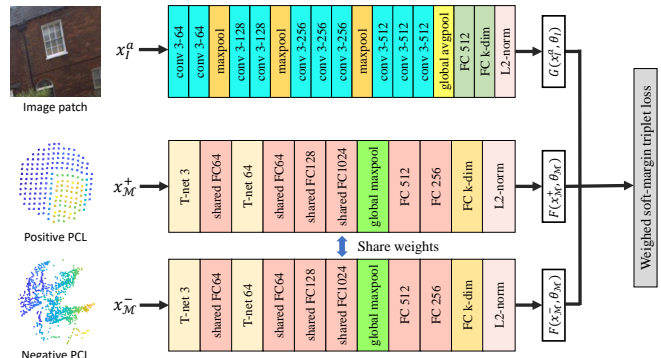


Fig. 1. Our triplet-like 2D3DMatch-Net.

directly from the 2D image and 3D point cloud. We first apply the SIFT detector on the query image I to extract a set of 2D keypoints $U = \{u_1, \dots, u_N \mid u_n \in \mathbb{R}^2\}$, and the ISS detector [8] on the 3D point cloud of the reference map M to extract a set of 3D keypoints $V = \{v_1, \dots, v_M \mid v_m \in \mathbb{R}^3\}$. Here, N and M are the total number of 2D and 3D keypoints extracted from the image I and point cloud M , respectively. Given the set of 2D image patches centered around each 2D keypoint and 3D local point cloud volume centered around each 3D keypoint, our 2D3D-MatchNet learns the corresponding set of 2D and 3D descriptors denoted as $P = \{p_1, \dots, p_N \mid p_n \in \mathbb{R}^D\}$ and $Q = \{q_1, \dots, q_M \mid q_m \in \mathbb{R}^D\}$ for each corresponding 2D and 3D keypoint in U and V . D is the dimension of the descriptor. Furthermore, the descriptors P and Q learned from our network yield a much smaller similarity distance $d(p, q)$ between a matching pair of 2D-3D descriptors (p, q) in comparison to the similarity distance $d(\bar{p}, \bar{q})$ between a non-matching pair of 2D-3D descriptors (\bar{p}, \bar{q}) , i.e. $d(p, q) \ll d(\bar{p}, \bar{q})$, thus establishing the 2D-3D correspondences between P and Q . Finally, the 2D-3D correspondences found from our 2D3D-MatchNet are used to estimate the absolute pose of the camera using a PnP algorithm. We run the PnP algorithm within RANSAC [27] for robust estimation.

B. Our 2D3D-MatchNet: Network Architecture

Our 2D3D-MatchNet is a triplet-like deep network that jointly learns the similarity between a given pair of image patch and local point cloud volume. The network consists of three branches as illustrated in Fig. 1. One of the branches learns the descriptor for the 2D image keypoint and the other two branches with shared weights learn the descriptor for the 3D point cloud keypoint. The inputs to the network are (1) image patches centered on the 2D image keypoints, and (2) local volume of point cloud within a fixed radius sphere centered on the 3D keypoints. Details on keypoints extraction and the definitions of image patch and point cloud sphere are given in Sec. IV-B.

The image patches and local volume of point clouds are fed into the network during training as tuples of anchor image patch, and positive and negative local volume point cloud. We denote the training tuple as $\{x_I^a, x_M^+, x_M^-\}$. Given a set

of training tuples, our network learns the image descriptor function $G(x_I; \theta_I) : x_I \mapsto p$ that maps an input image patch x_I to its descriptor p , and the point cloud descriptor function $F(x_M; \theta_M) : x_M \mapsto q$ that maps an input local point cloud volume x_M to its descriptor q . θ_I and θ_M are the weights of the network learned during training. More specifically:

Image Descriptor Function We design $G(x_I; \theta_I)$ as a convolutional neural network followed by several fully connected layers to extract the descriptor vector of an image patch. We use the well-known VGG network [28] as the basis of our image descriptor function network. We make use of the first four convolution block (conv1 ~ conv4) of VGG16 to make the network fit better for image patch with small size. Global average pooling is applied on the feature maps from conv4. Compared to the max pooling layer which is widely used after convolutional layers, global average pooling has the advantages of reducing the number of parameters and avoiding over-fitting. Two fully connected layers are concatenated at the end of the network to achieve the desired output descriptor dimension. The output descriptor vector is L2-normalized before feeding into the loss function.

Point Cloud Descriptor Function We use the state-of-art PointNet [29] as our point cloud descriptor function $F(x_M; \theta_M)$ to extract the descriptor vector of a local point cloud volume. The dimension of the last fully connected layer is changed to fit our feature dimension. The softmax layer in the end is replaced by a L2-normalization.

Loss Function for Training Our network is trained using the Triplet loss, so that the similarity distance $d_{pos} = d(G(x_I^a; \theta_I), F(x_M^+; \theta_M))$ between the matching anchor x_I^a and positive x_M^+ pair is small and much lesser than the similarity distance $d_{neg} = d(G(x_I^a; \theta_I), F(x_M^-; \theta_M))$ between the non-matching anchor x_I^a and negative x_M^- pair, i.e. $d(G(x_I^a; \theta_I), F(x_M^+; \theta_M)) \ll d(G(x_I^a; \theta_I), F(x_M^-; \theta_M))$. Specifically, Our network is trained with the weighted soft-margin triplet loss [22]:

$$\mathcal{L} = \ln(1 + e^{\alpha d}), \quad (1)$$

where $d = d_{pos} - d_{neg}$. We use this loss because it allows the deep network to converge faster and increase the retrieval accuracy [22]. In contrast to the basic Triplet loss [26], [30], it also avoids the need of selecting an optimal margin. In our experiment, we set $\alpha = 5$. We use the Euclidean distance between two vectors as the similarity distance $d(., .)$ in this work. During inference, we check the distance $d(G(x_I; \theta_I), F(x_M; \theta_M))$ between the descriptors of a given pair of image patch x_I and local point cloud volume x_M . x_I and x_M are deemed matching pair if the distance is lesser than a threshold, and non-matching pair otherwise.

C. Pose Estimation

The pose of the camera is computed from the putative set of 2D-3D correspondences obtained from our 2D3DMatchNet. Specifically, we obtain the 2D keypoints of the 2D query image with the SIFT detector, and the 3D keypoints of the 3D point cloud with the ISS detector. We compute the 2D and 3D keypoint descriptors with our network from the image

patches and local point cloud volume extracted around the keypoints. The similarity distance is computed for every pair of 2D and 3D keypoints, and we find the top K closest 3D point cloud keypoints for every 2D image keypoint. Finally, we apply the EPnP algorithm [9] to estimate the camera pose with all the putative 2D-3D correspondences. The EPnP algorithm is ran within RANSAC for robust estimation to eliminate outliers.

IV. DATASET

In this section, we present the creation of our benchmark dataset – *Oxford 2D-3D Patches Dataset*. The dataset contains a total of 432,982 image patch to pointcloud pairs, which allows sufficient training and evaluation for the 2D-3D feature matching task.

A. The Oxford 2D-3D Patches Dataset

Our Oxford 2D-3D Patches dataset is created based on the Oxford RobotCar Dataset [31]. The Oxford RobotCar Dataset collects data from different kinds of sensors, including cameras, Lidar and GPS/INS, for over one year. We use the images from the two (left and right) Point Grey Grasshopper2 monocular cameras, the laser scans from the front SICK LMS-151 2D Lidar, and the GPS/INS data from the NovAtel SPAN-CPT ALIGN inertial and GPS navigation system. Ignoring the traversals collected with poor GPS, night and rain, we get 36 traversals for over a year with sufficiently challenging lighting, weather and traffic conditions. We synchronize the images from the left and right cameras, and 2D laser scans from the Lidar with the timestamps, and get their global poses using the GPS/INS data. We remove camera and Lidar frames with small motion.

To simplify point cloud processing and reduce the detrimental effects of GPS jumps over long distance, we split each traversal into disjoint submaps at every 60m interval. Each submap contains the corresponding sets of left and right cameras and Lidar frames. A visualization of the reconstructed point cloud map from Lidar scans is illustrated in Fig. 2.

B. Training Data Generation

Keypoint Detection We build a point cloud based reference map from the laser scans for every submap, where the coordinates of the first laser scan is used as the reference frame. We detect the ground plane and remove all points lying on it. This is because the flat ground plane is unlikely to contain any good 3D keypoint and descriptor. The ISS keypoint detector is applied on the remaining point cloud to extract all 3D keypoints. We apply the SIFT detector on every image to extract all 2D keypoints.

2D-3D Correspondences To establish the 2D-3D correspondences, we project each ISS keypoint to all images within its view and find the nearest neighbour of SIFT keypoint in each image. To increase the confidence of the correspondences, we require the distance of the projected nearest neighbour to be smaller than 3 pixels and each ISS point must have at least SIFT correspondences in three

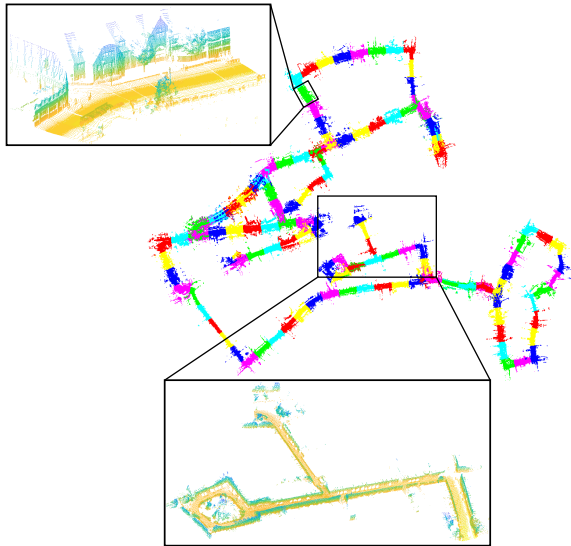


Fig. 2. Reconstructed point cloud map from Lidar. Different colors represent different submaps. Two zooming in examples of point cloud are shown, the top one indicates a 60m submap, the bottom shows the last unseen 10% of the path for testing.

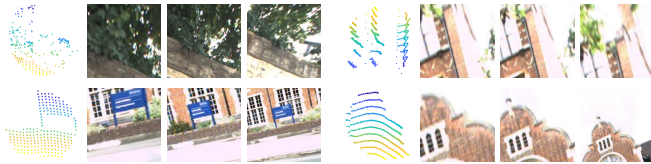


Fig. 3. Four examples of our dataset. The first image of each example shows the ISS volume. The other three are some corresponding SIFT patches across multiple frames with different scale, viewpoint and lighting.

different views within each submap. The ISS points and their corresponding SIFT points that satisfy these requirements are reserved for further processing.

ISS Volume and SIFT Patch Extraction We remove all ISS keypoints that are within 4m from a selected ISS keypoint in each submap, and remove all SIFT keypoints within 32 pixels from a selected SIFT keypoint in each image. We find all 3D points that are within 1m radius from each selected ISS keypoint, and discard those ISS keypoints with less than 100 neighboring 3D points. We discard SIFT keypoints with scale larger than a threshold value, since larger scale results in smaller patch size. In our experiments, we set this threshold value as 4 and the patch size at the basic scale as 256×256 . We vary the extracted patch size with the scale of the extracted SIFT keypoints for better scale invariance. In our experiments, we extract the ISS volume and its corresponding SIFT patch if the number of points within the ISS volume is larger than 100 and the SIFT patch is at suitable scale. We discard both the ISS volume and SIFT patch otherwise. Fig. 3 shows the visualization of several examples of the local ISS point cloud volumes and their corresponding image patches with different scales, viewpoints and lightings.

Data Pre-processing Before training, we rescale all the SIFT patches with different scales to the same size, i.e.

128×128 , and zero-center by mean subtraction. We subtract each point within each ISS point cloud volume with the associated ISS keypoint, thus achieving zero-center and unit norm sphere. Additionally, we pad the number of points to 1024 for each local volume in our experiments.

C. Testing Data Generation

Our objective during inference is to localize a query image based on the 2D-3D matching of the descriptors from the keypoints extracted from the image and point cloud. We test our trained network with reference submaps and images that are not used in training the network. We use the GPS/INS poses of the images as the ground truth pose for verification. The ground truth 2D-3D correspondences are computed as follows: (1) We detect all ISS keypoints from the point cloud of each submap and retain keypoints with more than 100 neighboring 3D points within 1m radius. (2) We detect SIFT keypoints on each image and extract the corresponding patches with scale smaller than the threshold value, i.e. 4 as mentioned above. (3) Each ISS keypoint is projected to all images within its view and the nearest SIFT keypoint with a distance smaller than 3 pixels is selected as the correspondence. We discard an ISS to SIFT keypoint correspondence if a nearest SIFT within 3 pixels is found in less than 3 image views.

V. EXPERIMENTS

In this section, we first outline the training process of our 2D3D-MatchNet that jointly learns both 2D image and 3D point cloud descriptors. Next, we describe the camera pose estimation given a query image. Then we evaluate our results and compare with different methods. Finally, we discuss and analyze the localization results of the proposed methods.

A. Network Training

Data splitting and evaluation metric As mentioned in Sec. IV-B, we split each traversal into a set of disjoint 60m submaps. We leave one full traversal for testing. For the remaining 35 traversals, we use the first 90% submaps of each traversal for training and leave the remaining 10% unseen for testing as in Fig. 2.

We evaluate the accuracy of the estimated pose by computing the position error and the rotational error from the ground truth poses. Similar to [32], we define the pose precision threshold as $(10m, 45^\circ)$. We measure the percentage of query images localized within this range and report the average position error and rotation error. We choose the threshold values to satisfy the high requirements for autonomous driving.

Network Training Our network is implemented in Tensorflow [33] with $2 \times$ Nvidia Titan X GPUs. We train the whole network in an end-to-end manner. For each triplet input, we choose an image patch as the anchor, and its corresponding 3D point cloud volume as positive sample. The negative point cloud volume is randomly sampled from the rest of the point clouds. We initialize the image descriptor network branch with VGG model pre-trained

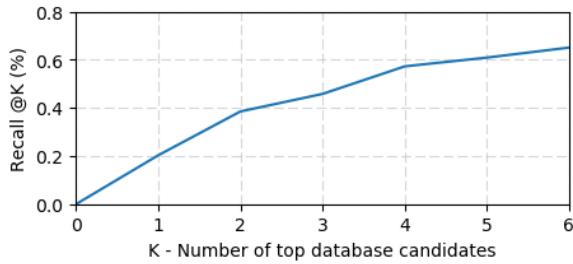


Fig. 4. The recall in top K .

on ImageNet [34]. Both descriptor extraction networks are optimized with Adam optimizer and the initial learning rate is 6×10^{-5} . The total training time is around two days.

We also explore the effect of the output feature dimension D on the performance of localization. In our experiments, we train and test with different D values, i.e. $D \in \{64, 128, 256\}$. The localization results from different descriptor dimensions are presented and discussed below.

B. Results

We test on one full traversal and the last unseen 10% from another six traversals. As mentioned in IV-C, we reconstruct the point cloud from GPS/INS data for each submap. We remove the redundant points from the ground plane. Next, we detect all the 3D keypoints and infer the corresponding descriptors from the point cloud descriptor network. All descriptors of the point cloud keypoints are stored in the database.

Network results Given a query image, we extract the 2D SIFT keypoints and feed all the corresponding image patches into the image descriptor network to get the descriptors of the query image. For each image descriptor, we find its top K nearest point descriptor from our database thus establishing the 2D-3D correspondences. In Fig. 4, it shows the recall from top-1 to top-6. The selection of K can largely effect the localization results. With a larger K , we have more point feature candidates for each image feature. Consequently, the RANSAC algorithm is more likely to find the correct match. On the other hand, a larger K unfavorably increases the number of iterations of RANSAC exponentially. Considering the trade-off, we choose $K = 5$ for our experiments.

Camera pose estimation Finally, we solve the camera pose using the EPnP algorithm [9]. For all images in each test submap, the localization results are presented in Tab. I. The unit of T and R error are meter and degree. The first result (@2014.06.26) reports the localization results on most submaps of the full run, except for those with bad point cloud maps due to GPS inaccuracy. We denote it as *full_test*. Others show the localization results on several different test submaps of each traversal across different times over one year. We denote them as *submap_test*. The ratio of successfully localized frames are shown in Fig. 6.

A qualitative visualization of the localization is presented in Fig.7. In these results, the output feature dimension D is set 128.

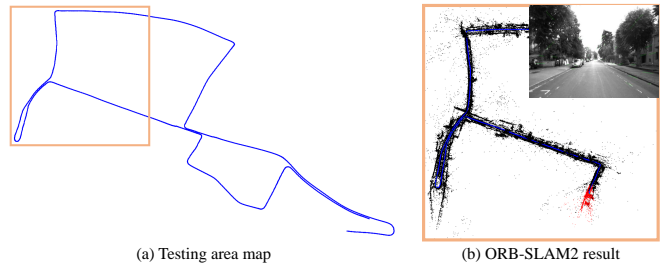


Fig. 5. The localization result on ORB-SLAM2. ORB-SLAM2 result: The red points are the failure position. An image of the failure case is shown.

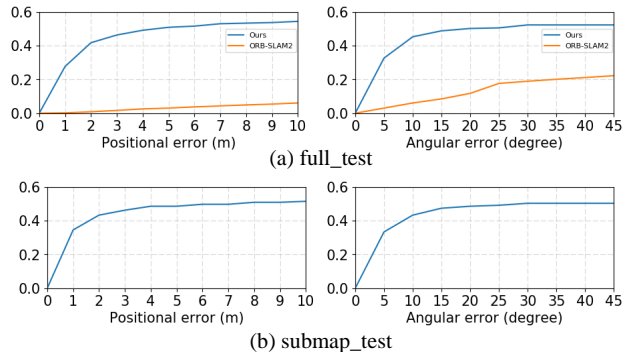


Fig. 6. The curves of ratio of successfully localized frames w.r.t thresholds. (a) *full_test*: results of testing on the *full_test* with both our method and ORB-SLAM2 [35]. (b) *submap_test*: average results of testing on the *submap_test*.

C. Evaluation and Comparison

For evaluation, we choose 8500 images collected in the overcast afternoon of July 14, 2015 with frame rate 16Hz. The path is around 3km with partial overlaps for loop closure detection, covering an spatial area of $845m \times 617m$. We test the performance of two well-known algorithms for the task of visual localization, i.e., ORB-SLAM2 [35] (traditional method) and PoseNet [15] (deep learning method).

We use ORB-SLAM2 algorithm to build the point cloud map of the testing area. However, it cannot succeed to build the whole area. Partial mapping result is shown in Fig. 5-(b). The localization result using the point cloud map by ORB-SLAM2 is shown in Fig. 6-(a). The large error is due to the inaccurate mapping. We observe that the map is unreliable when the images are captured near trees or at the turns. We train the PoseNet for visual localization. However, the localization error is huge. We argue that the PoseNet is not suitable in a large area with training data captured by cameras on a running vehicle, which do not provide rich variance on the angle and viewpoint. Furthermore, this method can not generalize to the unseen test sequences.

The existing algorithms fail to localize images within large-scale urban environments. In the next section, we show that our algorithm can successfully localize more than 40% of the images throughout the whole testing area.

TABLE I
THE LOCALIZATION RESULTS ON DIFFERENT TRAVERSALS FOR OVER ONE YEAR

Date and Time	2014.06.26 09:53:12	2014.07.14 15:16:36	2015.02.03 08:45:10	2015.04.24 08:15:07	2015.06.09 15:06:29	2015.07.14 16:17:39	2015.08.13 16:02:58
Test submaps	32	5	5	5	3	3	4
Test frames	7095	866	666	898	548	624	536
Average T error	1.41	1.34	1.88	1.67	1.68	1.44	1.71
Average R error	6.40	6.62	7.33	7.24	7.24	7.17	7.44

TABLE II
LOCALIZATION RESULTS ON DIFFERENT OUTPUT DESCRIPTOR
DIMENSION D ON 2015-02-13, 09:16:26

D	Test frames	Success frames	Average inliers	Average T error	Average R error
64	625	182	9	1.18	6.00
128		187	10	1.14	6.10
256		179	9	0.99	5.31

D. Analysis and Discussion

Generalization As can be seen in the Tab. I, the results of *submap_test* are slightly worse than the result of *full_test* since the testing area of *submap_test* is totally unseen. However, the results in unseen area are close to the results of seen area. It shows the good generalization of our proposed network.

Output Feature Dimension D We investigate the effect of the output feature dimension D . We test three submaps from another traversal on 2015, Feb 13, at 09:16:26. The localization result is shown in Tab. II. As we can see, the feature dimension D of 128 successfully localize more images than the other two. A higher feature dimension can better represent the image and point cloud, but on the other hand, it may also cause over-fitting since our patch size and point cloud volume are small. Considering the trade-off between accuracy and inference efficiency, we choose D as 128 in our experiments.

From the localization results, we show that our proposed method is able to estimate camera pose from a point cloud based reference map directly through 2D image to 3D point cloud descriptor matching using deep learning. There are two main cases where the localization is likely to fail. (1) The scene contains many trees: 3D points from trees are quite likely to be detected as key points due to their strong gradients, i.e. irregularity in shape. However, the SIFT keypoints on trees do not contain discriminative information. Consequently, wrong matches arose from patches and point clouds on trees. (2) The scene is dominated by flat building walls: buildings are always full of texture seen from image, and thus create many meaningful patches. However, points on smooth wall are less likely to be detected as keypoints. This leads to low 3D keypoints and descriptor candidates, which decreases localization performance.

VI. CONCLUSION

We presented a novel method for camera pose estimation given a 3D point cloud reference map of the outdoor environ-

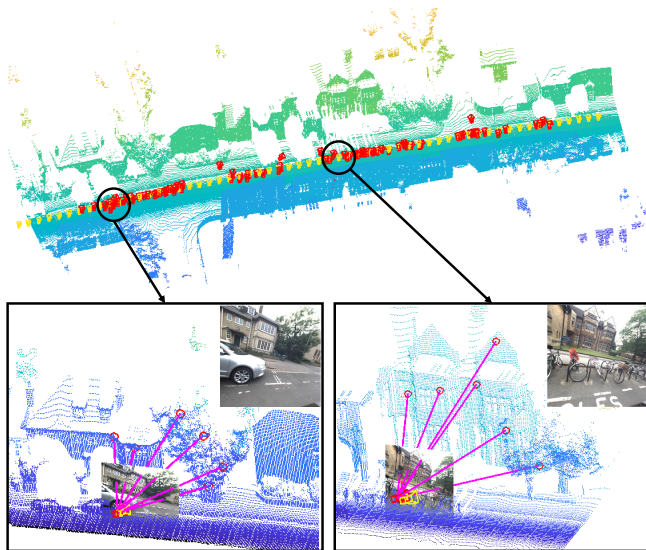


Fig. 7. A qualitative visualization of our camera pose estimation. Red camera: our estimated camera pose. Yellow camera: the ground-truth camera pose. Purple line: some predicted 2D-3D correspondences between 2D SIFT patches and 3D ISS volumes.

ment. Instead of the association of local image descriptors to points in the reference map, we proposed to jointly learn the image and point cloud descriptors directly through our deep network model, thus obtaining the 2D-3D correspondences and estimating the camera pose with the EPnP algorithm. We demonstrated that our network is able to map cross-domain inputs (i.e. image and point cloud) to a discriminative descriptor space where their similarity / dis-similarity can be easily identified. Our method achieved considerable localization results with average translation and rotation errors of 1.41m and 6.40 degree on the standard Oxford RobotCar dataset. In future work, we aim at an end-to-end network for camera pose estimation by incorporating the hand-crafted key point selection and the RANSAC algorithm into the network. Furthermore, we will enforce temporal consistencies on multiple continuous frames to help improve the localization accuracy.

VII. ACKNOWLEDGMENT

This research was supported in parts by the National Research Foundation (NRF) Singapore through the Singapore-MIT Alliance for Research and Technology's (FM IRG) research programme and a Singapore MOE Tier 1 grant R-252-000-637-112. We are grateful for the support.

REFERENCES

- [1] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [2] J. Engel, T. Schöps, and D. Cremers, "Lsd-slam: Large-scale direct monocular slam," in *European Conference on Computer Vision*, September 2014.
- [3] S. Agarwal, Y. Furukawa, N. Snavely, I. Simon, B. Curless, S. M. Seitz, and R. Szeliski, "Building rome in a day," *Communications of the ACM*, vol. 54, no. 10, pp. 105–112, 2011.
- [4] H. Lategahn, A. Geiger, and B. Kitt, "Visual slam for autonomous ground vehicles," in *The IEEE International Conference on Robotics and Automation*, May 2011.
- [5] M. Billinghurst, A. Clark, G. Lee, et al., "A survey of augmented reality," *Foundations and Trends® in Human-Computer Interaction*, vol. 8, no. 2-3, pp. 73–272, 2015.
- [6] D. Nistér, "A minimal solution to the generalized 3-point pose problem. on plane-based camera calibration: a general algorithm, singularities, applications," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2004.
- [7] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, 2004.
- [8] Y. Zhong, "Intrinsic shape signatures: A shape descriptor for 3d object recognition," in *THE IEEE International Conference on Computer Vision Workshops*, September 2009.
- [9] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Epnnp: An accurate o(n) solution to the pnp problem," *International Journal of Computer Vision*, vol. 81, no. 2, 2009.
- [10] T. Sattler, B. Leibe, and L. Kobbelt, "Efficient & effective prioritized matching for large-scale image-based localization," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 9, pp. 1744–1756, 2017.
- [11] Y. Li, N. Snavely, D. Huttenlocher, and P. Fua, "Worldwide pose estimation using 3D point clouds," in *European Conference on Computer Vision*, October 2012.
- [12] C. Valgren and A. J. Lilienthal, "Sift, surf & seasons: Appearance-based long-term localization in outdoor environments," *Robotics and Autonomous Systems*, vol. 58, no. 2, pp. 149–156, 2010.
- [13] I. Ulrich and I. Nourbakhsh, "Appearance-based place recognition for topological localization," in *The IEEE International Conference on Robotics and Automation*, April 2000.
- [14] F. Walch, C. Hazirbas, L. Leal-Taixe, T. Sattler, S. Hilsenbeck, and D. Cremers, "Image-based localization using lstms for structured feature correlation," in *The IEEE International Conference on Computer Vision*, October 2017.
- [15] A. Kendall, M. Grimes, and R. Cipolla, "Posenet: A convolutional network for real-time 6-dof camera relocalization," in *The IEEE International Conference on Computer Vision*, December 2015.
- [16] J. Shotton, B. Glocker, C. Zach, S. Izadi, A. Criminisi, and A. Fitzgibbon, "Scene coordinate regression forests for camera relocalization in rgb-d images," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2013.
- [17] J. Valentin, M. Nießner, J. Shotton, A. Fitzgibbon, S. Izadi, and P. H. Torr, "Exploiting uncertainty in regression forests for accurate camera relocalization," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2015.
- [18] E. Brachmann, A. Krull, S. Nowozin, J. Shotton, F. Michel, S. Gumhold, and C. Rother, "Dsac-differentiable ransac for camera localization," in *The IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [19] J. Schönberger, M. Pollefeys, A. Geiger, and T. Sattler, "Semantic visual localization," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2018.
- [20] S. Zagoruyko and N. Komodakis, "Learning to compare image patches via convolutional neural networks," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2015.
- [21] X. Han, T. Leung, Y. Jia, R. Sukthankar, and A. C. Berg, "Matchnet: Unifying feature and metric learning for patch-based matching," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2015.
- [22] S. Hu, M. Feng, R. M. H. Nguyen, and G. H. Lee, "Cvm-net: Cross-view matching network for image-based ground-to-aerial geolocalization," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2018.
- [23] W. Liao, M. Y. Yang, N. Zhan, and B. Rosenhahn, "Triplet-based deep similarity learning for person re-identification," *CoRR*, vol. abs/1802.03254, 2018.
- [24] N. N. Vo and J. Hays, "Localizing and orienting street views using overhead imagery," in *European Conference on Computer Vision*, October 2016.
- [25] Y. Guo, D. Tao, J. Yu, and Y. Li, "Deep similarity feature learning for person re-identification," in *Pacific-Rim Conference on Advances in Multimedia Information Processing*, September 2016.
- [26] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2015.
- [27] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," in *Readings in computer vision*. Elsevier, 1987, pp. 726–740.
- [28] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [29] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *The IEEE Conference on Computer Vision and Pattern Recognition*, July 2017.
- [30] G. Chechik, V. Sharma, U. Shalit, and S. Bengio, "Large scale online learning of image similarity through ranking," *The Journal of Machine Learning Research*, vol. 11, March 2010.
- [31] 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, 2017.
- [32] T. Sattler, W. Maddern, C. Toft, A. Torii, L. Hammarstrand, E. Stenborg, D. Safari, M. Okutomi, M. Pollefeys, J. Sivic, et al., "Benchmarking 6dof outdoor visual localization in changing conditions," in *Proc. CVPR*, vol. 1, 2018.
- [33] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al., "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *CoRR*, vol. abs/1603.04467, 2016.
- [34] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F.-F. Li, "Imagenet: A large-scale hierarchical image database," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2009.
- [35] 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.