

Single-Shot Global Localization via Graph-Theoretic Correspondence Matching

Shigemichi Matsuzaki¹, Kenji Koide², Shuji Oishi², Masashi Yokozuka², and Atsuhiko Banno²

Abstract—This paper describes a method of single-shot global localization based on graph-theoretic matching of instances between a query and the prior map. The proposed framework employs correspondence matching based on the maximum clique problem (MCP). The framework is potentially applicable to other map and/or query modalities thanks to the graph-based abstraction of the problem, while many of existing global localization methods rely on a query and the dataset in the same modality. We implement it with a semantically labeled 3D point cloud map, and a semantic segmentation image as a query. Leveraging the graph-theoretic framework, the proposed method realizes global localization exploiting only the map and the query. The method shows promising results on multiple large-scale simulated maps of urban scenes.

I. INTRODUCTION

Global localization is a problem to estimate a sensor pose in a prior map given a sensor observation or a sequence of observations, which is referred to as a *query*, and without prior information of its initial pose. It is a fundamental ability of mobile robots and autonomous vehicles for reliable and safe operation. For the task of global positioning, Global Navigation Satellite System (GNSS) is widely used. It is known, however, that the utility of GNSS is limited to open space without obstacles such as buildings, trees etc. that cause multi-path effect.

Various approaches to global localization have been proposed in a last few decades. In traditional Bayesian estimation such as Monte Carlo Localization (MCL) [1], a match between sensor readings and the prior map is evaluated through the observation model and incorporated as likelihood in iterative computation of the belief score on the robot pose. Single-shot methods such as global descriptor matching [2, 3] and local geometry matching [4, 5] instead calculate the most likely pose given a single query. In the former, a single global descriptor for the query is computed and used to retrieve from a database the most similar descriptor associated with a pose. The latter computes a set of descriptors that describe the local structure of the query, and find their correspondence with pre-computed feature points in the map to compute the pose via error minimization.

*This work was supported in part by a project commissioned by the New Energy and Industrial Technology Development Organization (NEDO). This work was also supported by JSPS KAKENHI (Grant Number 22K12214).

¹S. Matsuzaki is with Department of Computer Science and Engineering, Toyohashi University of Technology, Hibarigaoka 1-1, Tenpaku-cho, Toyohashi, Aichi, Japan matsuzaki@ais1.cs.tut.ac.jp

²K. Koide, S. Oishi, M. Yokozuka, and A. Banno are with the Department of Information Technology and Human Factors, the National Institute of Advanced Industrial Science and Technology, Umezono 1-1-1, Tsukuba, 3050061, Ibaraki, Japan k.koide@ieee.org

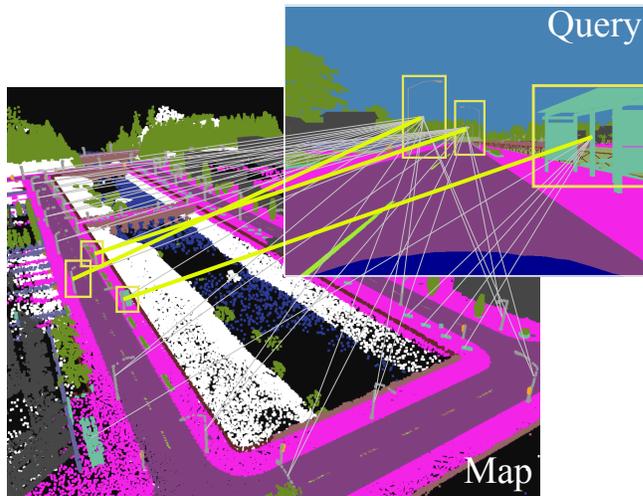


Fig. 1. Global localization is solved by estimating true correspondences between instances in a query and those in the prior map. The proposed graph-theoretic framework is potentially applicable to various query and map modalities.

Overall, the global localization methods essentially require matching between the query and the map. Regardless of the approaches, many of the existing global localization methods assume the query and the map in the same modality, i.e., the prior map must store data in the same representation as the query. This fact limits the applicability of the methods to a wide variety of different sensor and map modalities available. Recently, various map representations have emerged, such as point cloud maps, high-definition vector maps (HD maps), tagged maps like Google Map, etc. If those cross-modal data can be handled by a unified method in global localization, more flexible choices on the sensor and the map will be available for the global localization problem.

With such a prospect, we propose a framework of single-shot global localization, formulating it as instance correspondence matching between the query and the prior map and finding a pose in the map that best explains the correspondences. For the correspondence matching, we employ a graph-theoretic method that reduces the problem to the maximum clique problem (MCP). Specifically, we construct a graph called a *consistency graph* that represents local consistency between correspondence candidates, and find the most likely set of correspondences via MCP on the consistency graph, similarly to [6] etc. The motivation behind the formulation is to abstract the problem to a general graph

problem so that it is not constrained on a specific modality of the query and the map. The proposed method can potentially be applied to cross-modal global localization problems as long as a proper consistency criteria can be defined.

In the present work, we implement the framework using 3D semantic point cloud maps as a prior map and a semantic segmentation image as a query for correspondence matching. The consistency of the correspondences are evaluated based on the closeness in appearance from the pose calculated by a correspondence pair. Thanks to the graph-theoretic framework, global localization is realized with a combination of simple consistency criteria that does not require additional data for training. The method thus exploits only a 3D semantic map and a semantic segmentation query, in contrast to recent existing methods [5, 7] that depend on data-driven descriptors for correspondence matching. Our method can, therefore, be easily applied to various environments.

II. RELATED WORK

A. Global localization

There have been various approaches with different techniques and sensor modalities.

Iterative methods Traditionally, global localization has been solved in iterative algorithm such as Monte Carlo localization (MCL) [1] where the belief on the robot’s pose gradually converges by incorporating a temporal sequence of observations as the robot moves. This approach inherently requires motion of the robot, which limits the use case of the methods. In the present work, we opt for a single-shot global localization approach.

Global descriptor matching Global descriptor-based methods learn a discriminative descriptor per scene from the sensor observations. The pose is then estimated by matching the descriptor for the observation and pre-computed set of descriptors associated with a pose. This line of work includes Visual Place Recognition (VPR) [2, 3], geometry-based methods [8, 9], and hybrid methods [10, 11].

Local geometry matching In visual global localization, another popular approach is matching local visual features such as SIFT [12] and SURF [13] between a query image and a pre-computed map and optimizing the pose based on the feature correspondences [4, 14]. In 3D LiDAR-based geometry matching, local point descriptors such as [15] and point segments [5, 7] are often used.

Semantics-based localization Following the recent progress of semantic segmentation based on DNNs, global localization methods based on visual semantic information have also been actively studied. Results of semantic segmentation are reported to be more robust to appearance changes than low-level visual information such as feature points [16]. Semantic information has been, therefore, exploited to realize localization methods that are robust to appearance changes over time and viewpoint changes [17].

DNN-based pose regression DNNs are also exploited to directly estimate a sensor pose. Kendal et al. [18] first proposed to regress a 6 DOF pose directly from an image. Feng et al. [19] jointly trained the descriptors for 2D images

and 3D points to match data in those different modalities. This approach, again, is specialized to the data on which the network is trained and hard to apply to various scenes.

Problems in existing work and novelties of our work

Regardless of the approaches, most of the conventional methods use a query and a map in the same modality. On the other hand, the proposed method does not strictly assume the same modality. Our method, therefore, is potentially applicable to various types of maps such as point cloud maps and HD vector maps.

Although some recent works use DNNs to learn a strong global/local descriptor [2, 5], or to embed data from different modalities in a common feature space [19], such methods require a huge amount of training data. Moreover, pose regression using DNN is specialized to the data on which the model is trained, and difficult to apply to various scenes.

Our primary purpose is to develop a novel unified framework for global localization applicable to cross-modal settings. The pose is efficiently estimated using the graph-theoretic correspondence matching and simple criteria of appearance similarity. We show the effectiveness in an implementation using a semantic point cloud map and a semantic segmentation image.

B. Correspondence matching

Correspondence matching is a crucial task for some of the aforementioned methods. Random sample consensus (RANSAC) is traditionally used to extract correspondences out of candidates including outliers [14, 4]. Although it shows a high robustness against false positive matches, the number of iterations to achieve a reasonable result increases exponentially as a noise ratio increases due to its stochasticity, and thus is not applicable to cases of high outlier rates.

Graph-theoretic correspondence extraction has been recently shown effective under a high outlier rate where RANSAC fails to extract correspondences [6]. Bailey et al. [20] was one of the earliest to frame data association as MCP. The point association via MCP has been applied to global registration of point cloud [6] [21]. Koide et al. [22] proposed an automatic localization of fiducial tags on a 3D prior map using MCP-based data association.

To the best of our knowledge, this work is the first attempt to employ MCP-based correspondence matching for global localization. The proposed method leverages the graph-theoretic approach to efficiently find the most likely instance correspondences and carries out stable global localization.

III. GLOBAL LOCALIZATION VIA CORRESPONDENCE MATCHING BY MCP

The general algorithm of the proposed framework of global localization is shown in Algorithm 1. The core of the framework is to estimate the correspondences between the query instances and the map instances to calculate the pose that matches the correspondences the best. To this end, we employ graph-theoretic correspondence matching.

First, we generate candidates for correspondences between a query instance and a map instance. Here, we treat all

Algorithm 1 Global localization via graph-theoretic correspondence matching

Input: Map instance set $\mathcal{M} = \{m_i\}_{i=1}^{N_m}$, query instance set

$$\mathcal{Q} = \{q_j\}_{j=1}^{N_q}$$

Output: 3-DOF robot pose (x, y, θ)

- 1: Correspondence candidate set $\mathcal{C} = \mathcal{M} \times \mathcal{Q}$
 - 2: Vertex set $\mathcal{V} = \mathcal{C}$
 - 3: Edge set $\mathcal{E} = \emptyset$
 - 4: **for** every combination of $c_i, c_j \in \mathcal{C}$ **do**
 - 5: **if** `is_consistent`(c_i, c_j) **then**
 - 6: $\mathcal{E} \leftarrow \mathcal{E} \cup \{e_{ij}, e_{ji}\}$
 - 7: **end if**
 - 8: **end for**
 - 9: $\mathcal{C}^* \leftarrow \text{max_clique}(\mathcal{V}, \mathcal{E})$
 - 10: $(x, y, \theta) \leftarrow \text{calculate_pose}(\mathcal{C}^*)$
 - 11: **return** (x, y, θ)
-

possible pairs of an instance in the query and one in the map as valid candidates for correspondences.

Once the correspondence candidates are generated, we build a *consistency graph*, a node of which represents a candidate of correspondence between an instance in the query and one in the map, and an edge indicates consistency between two candidates represented by the nodes connected by the edge. All possible pairs of two correspondence candidates are evaluated based on certain consistency criteria, and if two candidates are evaluated as consistent, an edge is added between the nodes.

The correspondences are then estimated by solving MCP on the consistency graph, and the robot pose is calculated using the extracted correspondences so that it best satisfies the observations of the instances.

IV. IMPLEMENTATION ON 3D SEMANTIC MAP AND A SEMANTIC SEGMENTATION IMAGE

Here, we implement the proposed framework using 3D labeled point cloud map as a prior map, and a semantic segmentation image as a query. We consider estimation of a 3 DOF pose (x, y, θ) of the ground robot. We assume a robot with a camera facing forward and that its transformation relative to the robot's base is known, and a semantic segmentation image for an image taken by the camera is available. The overview of the workflow is shown in Fig. 2.

A. Clustering the map and the query instances

The instances of the 3D map are clustered and individually saved before the localization process. Here, the following four object classes are used for correspondence matching: *pole*, *traffic sign*, *traffic light*, *static*.

The map points are simply clustered based on their object class and spacial closeness. For the query image, the instances are clustered based on the connectivity. This may wrongly extract multiple instances overlapping in the image as one. We, however, treat such "connected" instances as an individual instance, expecting that they will be rejected in the subsequent process while there still remain sufficient

correspondences for pose estimation. Candidates of instance correspondences are all patterns of pairs of instances with the same object class between the map and the image.

B. Consistency evaluation

For every combination of two correspondence candidates, a pose is calculated so that it matches the size and the location of image instances and reprojection of the corresponding map instances. The consistency of the pair is then evaluated based on the similarity of the appearance between the query and the map reprojection from the calculated pose.

1) *Pose calculation using a correspondence pair:* Given two image instances, an angle α formed by two vectors to each image instance from the origin of the camera frame is determined. The region of the pose is then constrained to a circle C that has an inscribed angle α as shown in Fig. 3. Since there are two such circles, we determine C based on the positional relationship of the instances.

To determine the pose, we use the scale information of one of the correspondences. The relationship of the size in the 3D space s_m and that of its projection s_i is written as:

$$s_i = \frac{f \cdot s_m}{Z_c}, \quad (1)$$

where f denotes the focal length of the camera and Z_c denotes the Z coordinate of the instance in the camera frame. Since the sizes of both image instance and the map instance in a correspondence are known, Z_c can be estimated as:

$$Z_c = \frac{f \cdot s_m}{s_i}. \quad (2)$$

The set of points that satisfy eq. (2) also forms a circle with a radius $D = \sqrt{X_c^2 + Z_c^2}$. The camera pose can, therefore, be calculated as an intersection of circle C and a circle that centers at the location of the map instance and has radius D . In this way, at most two intersections are yielded. In such a case, we calculate the expected size of the projection of another map instance using eq. (1) from the two poses, and compare them with the size of the corresponding image instance. The pose that provides a smaller error of the size is selected as the solution.

2) *Consistency criteria:* The consistency is evaluated based on two criteria, *background similarity* and *instance proximity*. The latter is relevant to the spacial consistency, which is often employed in MCP-based correspondence matching such as [6] and [22]. In our problem, however, this criterion is not sufficient because there may be many wrong correspondences that happen to satisfy it. The background similarity criterion is, therefore, employed before the instance proximity to filter out unlikely correspondences based on the appearance similarity. Using the two criteria, we effectively evaluate consistency of the correspondences.

Background similarity The purpose of the background similarity criterion is not to strictly detect images similar to the query, but to quickly filter out obvious outliers. To this end, we employ a simple histogram descriptor.

Fig. 4(a) shows an overview. First, the bottom of the image is cut off as this part does not convey much information.

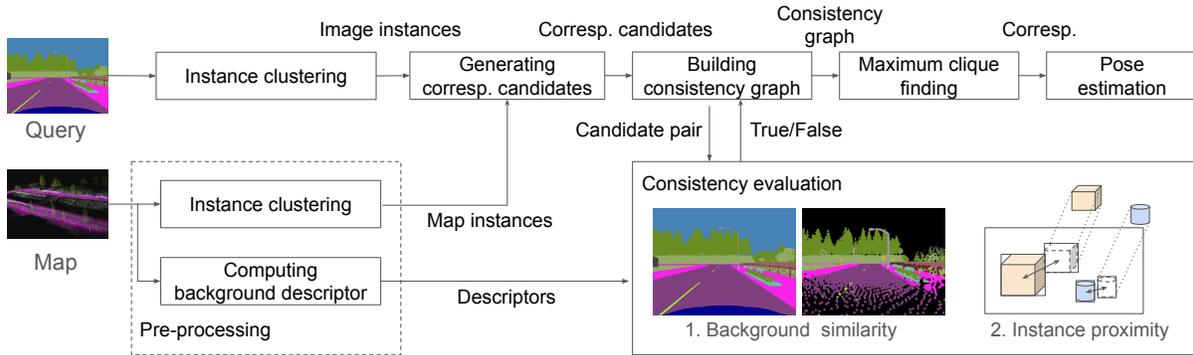


Fig. 2. Overview of the proposed framework. As pre-processing, the map instances are clustered, and the grid histogram descriptors (see IV-B.2) are computed, and stored as a database. Given a query image, instances are clustered and correspondence candidates are generated as all possible combinations of a query and a map instances. For every pair of the candidates, a pose is calculated (IV-B.1) and the consistency is evaluated based on two criteria 1) background similarity, and 2) instance proximity to construct a consistency graph (IV-B.2). The most likely set of correspondences is yielded by solving MCP on the consistency graph, and the pose is calculated based on it (IV-C).

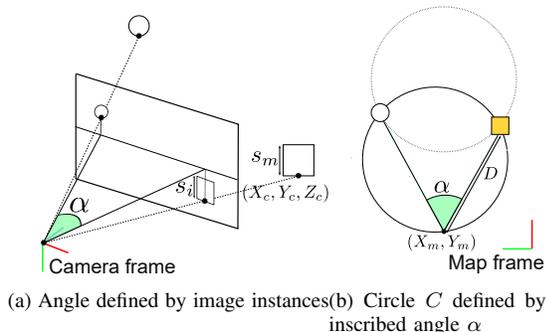


Fig. 3. The candidate regions of the robot pose based on the inscribed angle and the size of the instance. (b) The translation of the camera (X_m, Y_m) in the map frame is determined as a point on C where the distance to an object is $D = \sqrt{X_c^2 + Z_c^2}$.

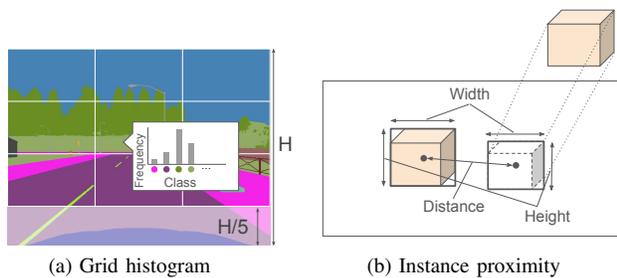


Fig. 4. Consistency criteria. (a) The grid histogram descriptor is a normalized vector of the local histograms of object classes. (b) The instance proximity evaluates a spatial relationship between the bounding boxes around the image instance and the projection of the map instance.

The remaining part of the image is divided into grids with a fixed size. For each grid, a histogram of the object classes is generated and L2-normalized. We chose nine background object classes: *building*, *fence*, *road line*, *road*, *sidewalk*, *vegetation*, *wall*, *bridge*, *terrain*. The normalized histograms from the grids are concatenated into $(h \times w \times 9)$ -dimensional vector, where h and w denotes the number of grid in the vertical and the horizontal axis, respectively. The vector is

then L2-normalized again to form a descriptor. Similarity of two descriptors are given as a dot product, or equivalently a cosine similarity with a value range of $[0, 1]$.

For fast evaluation, we pre-compute the descriptors for the map by projecting the 3D map from the poses uniformly sampled in the map.

Instance proximity The instance proximity is evaluated by the distance and the size differences of bounding boxes around the projected map instance and the image instance (see Fig. 4(b)). The distance between them are defined as the Euclidean distance between the centroids of the bounding boxes. The size differences are calculated with height and width lengths of the bounding boxes. If the distance, the height and the width difference are all less than a threshold for both of the correspondences in the given pair, the pair of correspondence candidates is considered consistent.

Here, we set the threshold adaptively to the Z coordinate value Z_c of each map instance in the camera frame. We empirically define the distance threshold θ_d and the size threshold θ_s as follows:

$$\begin{aligned} \theta_d &= \min(m_d, \max(m_{min}, dm_d/Z_c)) \\ \theta_s &= \min(m_s, \max(m_{min}, dm_s/Z_c)), \end{aligned} \quad (3)$$

where m_d and m_s denote the maximum possible threshold of the instance distance and the size in pixel, respectively. The idea behind these definitions is to set the thresholds of objects at d [m] from the sensor to m_d [px] and m_s [px], and decrease it as Z_c becomes larger. The thresholds are clipped with the minimum value of m_{min} [px] and the maximum value of m_d and m_s . Here, we set d to 8 and m_{min} to 5.

C. Pose estimation

A consistency graph is constructed based on the aforementioned consistency criteria, and the most likely set of correspondences are yielded by solving as described in IV. The final pose is then calculated by minimizing the reprojection error of representation points of the correspondences. A representation point of an image instance is a point that has the largest y value in the image coordinate, and that of

a map instance is the point with the smallest z value in the map coordinate, representing their bottom point. Although this way does not provide exact matching of the parts of the instances, it does not affect much the calculation of the camera pose. An initial pose is given as an average of poses calculated by the correspondences in the clique during the phase of building the consistency graph (see IV-B.1). The Levenberg-Marquardt algorithm [23, 24] with Huber kernel [25] is applied to the pose optimization.

D. Pose verification via correspondence check

Ideally, the maximum clique of the consistency graph represents a set of correct correspondences. In practice, however, we found that it is often not the case due to rejection of the correct correspondences mainly caused by error in pose calculation, while evaluating wrong ones as consistent.

To mitigate the problem, we introduce a verification step. After constructing the consistency graph, we calculate multiple correspondence set by iteratively applying MCP solver, calculating a pose with each set, and removing edges between the vertices corresponding to the members of the found maximum clique. To verify each correspondence set, we count the number of image instances that have a map instance projected close to it from the calculated pose based on the proximity criterion. The correspondence set with the maximum count is treated as the result. Top N results are yielded by taking results with five largest counts.

V. EXPERIMENTS

A. Experimental setup

The experiments were conducted on a laptop with an Intel Core i7-10750H (12 threads) and 16 GB of RAM.

We used simulated point cloud maps of Town01 to 05 from CARLA simulator [26] for evaluation. The description of the environments is shown in Table I. We created point cloud maps using the scripts shared by the authors of [27]. For each map, descriptors are pre-computed from the poses uniformly sampled with a stride of 2 m in the x and the y axes of the map frame and 30° in the yaw angle.

The samples of query were generated from the spawn points of vehicles defined in CARLA. A camera was mounted on the simulated vehicle at the height of approx. 1.5 m from the robot’s base. Since the proposed method estimates a 3 DOF pose and assumes fixed height of the camera, we choose the poses with the height between 1.46 m and 1.90 m as the samples. We use ground truth semantic segmentation images generated by CARLA.

The thresholds of the background similarity (BS), the instance distance and the instance size in the instance proximity (IP) are empirically set to 0.9, 110, 50, respectively.

B. Baselines

As baselines, we use brute-force descriptor search (BF) and RANSAC-based correspondence matching. The former searches for the pre-computed descriptor that has the highest

¹The descriptions are from https://carla.readthedocs.io/en/latest/core_map/

TABLE I
DESCRIPTORS OF THE CARLA TOWNS¹

Town	Summary
01	A basic town layout consisting of "T junctions".
02	Similar to Town01, but smaller.
03	The most complex town, with a 5-lane junction, a roundabout, unevenness, a tunnel, and more.
04	An infinite loop with a highway and a small town.
05	Squared-grid town with cross junctions and a bridge. It has multiple lanes per direction.

similarity score with the query and returns the discretized pose that the descriptor is associated with. In the latter, RANSAC is used for correspondence finding instead of the proposed graph-theoretic approach. For a randomly sampled pair of correspondences, a pose is calculated as described in IV-B.1 and the number of the correspondences consistent with the pose is counted. We use the same background similarity threshold as the proposed method to reject the pair, and the instance proximity threshold to find the correspondences. The number of iteration is set to 50000.

C. Evaluation conditions

We evaluate the localization performance using three conditions: ± 5 , ± 10 , and *front drift*. For all of the conditions, the threshold of yaw error is $\pm 30^\circ$. ± 5 denotes the case where the lateral and longitudinal errors are within ± 5 m. Likewise, ± 10 means within ± 10 m of lateral and longitudinal errors. *Front drift* denotes the case where the error in the x axis of the robot’s frame is within ± 200 m and ± 5 m in the y axis. This is typically the case when the estimated pose is on the same road as the ground truth and facing the same direction, but is drifted in the longitudinal direction in the robot frame.

D. Localization results

Table II shows the results of global localization. The values indicate the number of samples that resulted in fulfilling the evaluation conditions in Top 1, Top 3, and Top 5 results calculated as described in IV-D. In Town01 and 02, Top 5 estimation for approximately 65 to 70 % of the samples fulfilled ± 5 , and more than 80 % of them fulfilled *front drift*. In Town 03 and 05, Top 5 results of approximately 40 % samples fulfilled ± 5 . The result of Town04 was significantly worse than the other maps, with 18.5 % of ± 5 and 30.1 % of *front drift* in Top 5 results, respectively. We discuss the cause of the low performance in V-G.

E. Comparison with the baselines

Table III shows the comparative results with the baselines. The values are the number of Top 1 results of ± 5 and ± 10 . The proposed method surpassed the baseline methods in all conditions but ± 10 of BF in Town05. It is worth noting that the brute-force search did not result in accurate pose estimation in most of the environments. Nevertheless, combined with our framework, it resulted in better performance. RANSAC was also resulted in worse results than the proposed method. This is an expected result as most of

TABLE II
RESULTS OF GLOBAL LOCALIZATION

Town		01	02	03	04	05
Samples		254	101	197	259	234
± 5	T1	133 (52.4%)	60 (59.4%)	63 (32.0%)	36 (13.7%)	76 (32.5%)
	T3	157 (61.8%)	68 (67.3%)	73 (37.1%)	45 (17.4%)	93 (39.7%)
	T5	167 (65.7%)	71 (70.3%)	79 (40.1%)	48 (18.5%)	98 (41.9%)
± 10	T1	142 (55.9%)	65 (64.4%)	72 (36.5%)	49 (18.9%)	81 (34.6%)
	T3	170 (66.9%)	74 (73.3%)	85 (43.1%)	62 (23.9%)	100 (42.7%)
	T5	183 (72.0%)	77 (76.2%)	94 (47.7%)	67 (25.9%)	106 (45.3%)
FD	T1	175 (68.9%)	75 (74.3%)	95 (48.2%)	53 (20.5%)	87 (37.2%)
	T3	196 (77.2%)	81 (80.2%)	104 (52.8%)	76 (27.8%)	110 (47.0%)
	T5	207 (81.5%)	86 (85.1%)	114 (57.9%)	82 (30.1%)	123 (52.6%)

* T1, T3, and T5 denote Top 1, Top 3, and Top 5, respectively. *FD* denotes *Front drift*. **Green** : > 70 %, **Red** : < 30 %

TABLE III
COMPARISON WITH THE BASELINES

	Town no.	01	02	03	04	05
	Samples	254	101	197	259	234
BF	± 5	99	35	49	32	70
	± 10	107	34	62	47	98
RANSAC	± 5	17	7	5	2	6
	± 10	22	10	17	11	12
Proposed	± 5	133	60	63	36	76
	± 10	142	65	72	49	81

the correspondence candidates are outliers in the tasks, and RANSAC struggles with finding a consistent pair via random sampling. In contrast, the proposed method efficiently found the correspondences leveraging the graph-theoretic approach.

F. Ablation studies on the consistency criteria

Table IV shows the results of an ablation study on the consistency criteria to examine their contribution. To test IP for Town02 to 05, we loosened BS threshold to 0.5 (denoted as IP+BS(0.5)) instead of using only IP, because we found disabling BS significantly slowed down the computation speed due to increased number of transformation of map

TABLE IV
ABLATION ON THE CONSISTENCY CRITERIA

Town no.	01	02	03	04	05
Samples	254	101	197	259	234
BS(0.9)	40	18	6	1	8
IP/IP+BS(0.5)	18*	8	4	11	15
IP+BS(0.9)	133	60	63	36	76

* Only IP

instances in IP computation.

As a result, the performance was significantly worse in BS and IP (Town01) / IP+BS(0.5) (02-05). The combination of those simple consistency criteria resulted the best to effectively find the correspondences and calculate the pose.

G. Discussion

We saw promising results of the proposed framework especially in Town01 and 02. Notably, the proposed method requires only the semantic 3D map and semantic segmentation query. Although the building blocks for consistency evaluation are simple, our system provided reasonable results. We suppose this is thanks to the graph-theoretic correspondence matching approach which is capable of effectively extracting correspondences under high outlier rate.

Limitations Despite the powerful graph-theoretic approach, there were also many failure cases. They mainly stem from the error in the pose estimation using a pair of correspondences. The current pose calculation is sensitive to noise in instance clustering because it exploits the relationship of the scales of instances to estimate the distance from the robot pose to the map instances, and assumes that both the image instance and the map instance are complete. In reality, however, the assumption often does not hold. This problem will also affect real-world application because, in practice, the query semantic segmentation will be provided by an estimator based on DNNs etc., which inevitably has estimation noise. To deal with it, we should introduce pose calculation that is more robust to such noise.

The causes of the low accuracy for Town04 are twofold. First, the number of image instances are fewer than other Towns. Due to the unstable pose calculation, a few image instances in an image are generally rejected in MCP, and it is thus more possible to result in no correct correspondences detected when an image has few instances. Second, the mountains and terrains seen far away in the query were not mapped due to the limited range of the point cloud map, and the descriptors of the query and the map had a large difference in many samples. This effect was especially significant in Town04 consisting of open highway scenes with mountains. Although the simple descriptor provided reasonable results in tandem with the proposed framework, we should reconsider its design to mitigate such a problem.

VI. CONCLUSION

We proposed a framework for global localization based on graph-theoretic matching of instance correspondences between a query and a prior map. We implemented the framework on semantically labeled point cloud maps and semantic segmentation images as queries. Leveraging the efficient graph-theoretic correspondence matching with the simple consistency criteria, we presented promising preliminary results on the global localization task in simulated urban environments. This research direction will lead to localization ability of robots with more flexibility of the map and sensor modalities.

REFERENCES

- [1] F. Dellaert, W. Burgard, D. Fox, and S. Thrun, "Using the CONDENSATION algorithm for robust, vision-based mobile robot localization." In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, IEEE, 1999, pp. 588–594.
- [2] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN Architecture for Weakly Supervised Place Recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 6, pp. 1437–1451, 2018.
- [3] A. Khaliq, M. Milford, and S. Garg, "MultiRes-NetVLAD: Augmenting Place Recognition Training with Low-Resolution Imagery," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3882–3889, 2022.
- [4] M. Geppert, P. Liu, Z. Cui, M. Pollefeys, and T. Sattler, "Efficient 2D-3D matching for multi-camera visual localization." In *Proc. of the IEEE International Conference on Robotics and Automation*, 2019, pp. 5972–5978.
- [5] R. Dubé, A. Cramariuc, D. Dugas, H. Sommer, M. Dymczyk, J. Nieto, R. Siegwart, and C. Cadena, "SegMap: Segment-based mapping and localization using data-driven descriptors," *International Journal of Robotics Research*, vol. 39, no. 2-3, pp. 339–355, 2020.
- [6] H. Yang, J. Shi, and L. Carlone, "TEASER: Fast and Certifiable Point Cloud Registration," *IEEE Transactions on Robotics*, vol. 37, no. 2, pp. 314–333, 2021.
- [7] A. Cramariuc, F. Tschopp, N. Alatur, S. Benz, T. Falck, M. Bruhlmeier, B. Hahn, J. Nieto, and R. Siegwart, "SemSegMap – 3D Segment-based Semantic Localization." In *Proc. of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, IEEE, 2021, pp. 1183–1190.
- [8] H. Yin, L. Tang, X. Ding, Y. Wang, and R. Xiong, "LocNet: Global Localization in 3D Point Clouds for Mobile Vehicles." In *Proc. of the IEEE Intelligent Vehicles Symposium*, IEEE, 2018, pp. 728–733.
- [9] K. P. Cop, P. V. K. Borges, and R. Dube, "Delight: An Efficient Descriptor for Global Localisation Using LiDAR Intensities." In *Proc. of the IEEE International Conference on Robotics and Automation*, IEEE, 2018, pp. 3653–3660.
- [10] J. L. Schonberger, M. Pollefeys, A. Geiger, and T. Sattler, "Semantic Visual Localization." In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, IEEE, 2018, pp. 6896–6906.
- [11] S. Ratz, M. Dymczyk, R. Siegwart, and R. Dube, "OneShot Global Localization: Instant LiDAR-Visual Pose Estimation." In *Proc. of the IEEE International Conference on Robotics and Automation*, 2020, pp. 5415–5421.
- [12] D. G. Lowe, "Distinctive Image Features from Scale-Invariant Keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [13] H. Bay, T. Tuytelaars, and L. Van Gool, "SURF: Speeded Up Robust Features." In *Proc. of the European Conference on Computer Vision*, 2006, pp. 404–417.
- [14] S. Se, D. G. Lowe, and J. J. Little, "Vision-based global localization and mapping for mobile robots," *IEEE Transactions on Robotics*, vol. 21, no. 3, pp. 364–375, 2005.
- [15] R. B. Rusu, N. Blodow, and M. Beetz, "Fast Point Feature Histograms (FPFH) for 3D registration." In *Proc. of the IEEE International Conference on Robotics and Automation*, IEEE, 2009, pp. 3212–3217.
- [16] S. Orhan, J. J. Guerrero, and Y. Bastanlar, "Semantic Pose Verification for Outdoor Visual Localization with Self-supervised Contrastive Learning." In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, IEEE, 2022, pp. 3989–3998.
- [17] X. Guo, J. Hu, J. Chen, F. Deng, and T. L. Lam, "Semantic Histogram Based Graph Matching for Real-Time Multi-Robot Global Localization in Large Scale Environment," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 8349–8356, 2021.
- [18] A. Kendall, M. Grimes, and R. Cipolla, "PoseNet: A Convolutional Network for Real-Time 6-DOF Camera Relocalization." In *Proc. of the IEEE International Conference on Computer Vision*, IEEE, 2015, pp. 2938–2946.
- [19] M. Feng, S. Hu, M. H. Ang, and G. H. Lee, "2D3D-Matchnet: Learning To Match Keypoints Across 2D Image And 3D Point Cloud." In *Proc. of the IEEE International Conference on Robotics and Automation*, vol. 2019-May, IEEE, 2019, pp. 4790–4796.
- [20] T. Bailey, E. Nebot, J. Rosenblatt, and H. Durrant-Whyte, "Data association for mobile robot navigation: a graph theoretic approach." In *Proc. of the IEEE International Conference on Robotics and Automation*, IEEE, 2000, pp. 2512–2517.
- [21] J. Shi, H. Yang, and L. Carlone, "ROBIN: a Graph-Theoretic Approach to Reject Outliers in Robust Estimation using Invariants." In *Proc. of the IEEE International Conference on Robotics and Automation*, IEEE, 2021, pp. 13 820–13 827.
- [22] K. Koide, S. Oishi, M. Yokozuka, and A. Banno, "Scalable Fiducial Tag Localization on a 3D Prior Map via Graph-Theoretic Global Tag-Map Registration." In *Proc. of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, IEEE, 2022.
- [23] K. Levenberg, "A method for the solution of certain non-linear problems in least squares," *Quarterly of Applied Mathematics*, vol. 2, no. 2, pp. 164–168, 1944.
- [24] D. W. Marquardt, "An Algorithm for Least-Squares Estimation of Nonlinear Parameters," *Journal of the Society for Industrial and Applied Mathematics*, vol. 11, no. 2, pp. 431–441, 1963.
- [25] P. J. Huber, "Robust Estimation of a Location Parameter," *The Annals of Mathematical Statistics*, vol. 35, no. 1, pp. 73–101, 1964.
- [26] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An Open Urban Driving Simulator." In *Proc. of the Conference on Robot Learning*, PMLR, 2017, pp. 1–16.
- [27] J. E. Deschaud, D. Duque, J. P. Richa, S. Velasco-Forero, B. Marcotegui, and F. Goulette, "Paris-carla-3d: A real and synthetic outdoor point cloud dataset for challenging tasks in 3d mapping," *Remote Sensing*, vol. 13, no. 22, pp. 1–24, 2021.