

CMPUT 631: Project Proposal

Performance Enhancement of RTAB-map Utilizing ORB-SLAM2 Depth Information

Islam Ali
iaali@ualberta.ca

Zifei Jiang
zifei.jiang@ualberta.ca

Abstract

SLAM (Simultaneous Localization and Mapping) became one of the main building blocks of any modern robotic system due to of the inevitable need for robust and reliable localization and mapping engine in the robotic perception subsystem. A number of trials have been done during the last few decades to provide reliable solutions but each suffered a downside in a certain domain or under certain environmental conditions. Two famous SLAM algorithms (namely RTAB-map and ORB-SLAM2) are considered the industry standard as they provide state-of-the-art performance. However, each suffer critical performance issues in some testing scenarios. In this work, we explore the possibility of enhancing the performance of RTAB-map by utilizing the depth information of ORB-SLAM2. The proposal starts by defining the problem in hand and by giving a quick background and literature review of the efforts exerted in this track. Then, the methodology and procedure are discussed thoroughly to provide a road map of the project activities. Finally, the proposal conclude by providing a mechanism for performance evaluation as well as discussing the significance of this work.

1 Problem Definition

Recently, modern autonomous robotic systems are playing a vital role in a wide spectrum of industrial and non-industrial applications. One problem that is related to such application is robot perception of the environment, as such how it models the environment and how it localize itself in it accurately and reliably. To answer these questions, SLAM was introduced a few decades ago as a method for simultaneous localization and mapping. Major advancements in this field was introduced by the introduction of two state-of-the-art algorithms which are *RTAB-map* [1] and *ORB-SLAM* with its two versions [2][3]. According to a number of comparative studies [4][5], both algorithms were shown to outperform both classical and modern SLAM algorithms in some essence. However, these studies also pointed out performance issues in both under some conditions related to the environment and the dynamics of the camera/sensor motion. Having a deeper look at both algorithms and the their points principle of operation, one can infer some complementary features of both, and can deduce a lot of opportunities for integration. Interestingly, it was shown in [1] that the usage of ORB-SLAM2 as an odometry method for RTAB-map can result in an enhanced performance. In this work, we explore the possibility of embedding the optimized depth information from ORB-SLAM2 into RTAB-map to enhance its performance.

2 Background

SLAM (Simultaneous Localization and Mapping) is the process of generating a model of the environment (a map) and localizing the camera inside it [6]. By localization we refer to the estimation of the 6DoF of the moving camera (orientation and position). Many sensors are used in SLAM such as LiDARs, Cameras, and Radars [7]. In this work, we focus on the usage of vision-based SLAM and possibly its integration with Laser-based SLAM.

The conventional pipeline of SLAM consists of two major stages, the first one is responsible for sensor abstraction and processing, which includes features extraction and tracking as well as any long-term data association steps such as bundle adjustment or loop closure [7][6]. The second step is to use the associated data to estimate a map or a model of the environment and, at the same time, estimate the 6DoFs (complete pose) of the camera inside the generated map reliably.

3 Literature Review

3.1 Visual SLAM

A wide spectrum of research has been conducted into the SLAM problem with the objective of enhancing its performance, allow for long-term operation, and enable portable and limited resources applications to utilize it in real-time. SLAM algorithms can be categorized to either feature-based method and direct methods [8]. Feature-based methods depend on extracting features and tracking them through different captured images. The estimation of the camera motion and the map is then determined using either filter-based methods (such as EKF in MonoSLAM) or using bundle adjustment optimization (such as BA used in PTAM). It was shown that BA method can provide better performance due to its ability to use more feature key points when compared to the complex EKF-based solution [9].

On the contrary, direct feature-less methods are usually used on images directly without any abstraction [8] in order to achieve real-time operation and eliminating the accumulated error usually occurring in the aforementioned feature-based methods. Direct methods differ, mainly, in the density of the map and can be categorized based on that to either dense, semi-dense, or sparse methods [8]. Methods such as DTAM and LSD-SLAM are examples of the trials done to achieve direct SLAM by means of dense and semi-dense maps while other methods such as SVO and DVO are more advanced steps towards fully direct SLAM algorithms utilizing sparse maps [6][8].

3.2 RTAB-Map

RTAB-Map, (Real-Time Appearance-Based Mapping), is an open source Graph-Based SLAM system. RTAB-Map recently has been used as a platform to comparing various of SLAM algorithms due to its deep integration with ROS [1]. Based on ROS, developers are able to modify and extend the system easily. In the mean time, RTAB-Map provide two classic visual odometry algorithms Frame-To-Map(F2M) and Frame-To-Frame(F2F) [10]. F2F method registers the new frame against last frame while F2M registers the new frame against a local map features created from past key frames. The mapping process of RTAB-Map is

separated from odometry part and to reduce the calculation burden. RTAB-Map create Point Cloud using the depth image directly without any optimization.

3.3 ORB-SLAM 1 & 2

ORB-SLAM 1 is a monocular SLAM system proposed in 2015 [2]. It use ORB features for all nodes in the SLAM system such as tracking, mapping, relocalization and loop closing. ORB features provide high calculation efficiency to enable the real time operation in large environment. Based on ORB-SLAM 1, ORB-SLAM 2 is developed to work with RGBD and stereo camera. ORB-SLAM 2 has three main parallel threads: tracking, local mapping and loop closing. Since ORB-SLAM 2 optimize both state of the camera and key points in the map, it achieves the state-of-art in accuracy in SLAM systems.

3.4 RTAB-Map vs. ORB-SLAM 1&2

A number of trials have been made to evaluate the performance of SLAM algorithms (specially RTAB-map and ORB-SLAM 1&2) utilizing unified test procedures. These trials differ in the test environment (indoor vs. outdoor), the sensors used (monocular vs. stereo vs. RGB-D), and the test scenario used (testbeds vs. datasets vs. synthetic environment). In this subsection, we summarize the results of these studies.

ORB-SLAM exhibits correct estimation of the trajectory in terms of the traveled distance. However, it seems to lose track sometimes in case of doing turns [4] or when the camera is moving too fast. When ORB-SLAM loses track, it takes too long to re-initialize and re-start [11] the estimation procedure, resulting in degraded performance from the *recovery after failure* perspective.

ORB-SLAM2 was found to outperform RTAB-map in terms of odometry measurements [11] when tested with both stereo and RGB-D sensors. From the mapping perspective, it was shown in [5] that although the small number of map points and the little details provided, the method produced a smaller number of outliers. Also, ORB-SLAM2 was proven to work better in indoor environments [11][5]. In [12], it was observed that ORB-SLAM2 loses potential loop closures and thus the trajectory drift becomes visible. Additionally, it re-initialize several times during operation due to the number of tracked features going under a certain threshold.

RTAB-Map was found to provide better trajectory estimation but not distance measurement. It was also found to provide better map estimation and perform better in terms of loop-closure detection [11], and have a small number of outliers due to its rigorous filtering mechanism [13]. This is due to the dense map utilized in RTAB-Map and the BoW-based method used for loop closure. In [13], it was observed that RTAB-map will generate repeated surfaces due to its odometry noise. However, the performance is enhanced when using ORB-SLAM2 as the odometry method for RTAB-Map [1]. RTAB-Map showed superior performance in terms of RMS error but was outperformed by ORB-SLAM2 in terms of the max. trajectory error [5]. Finally, it was reported in [4] that RTAB-Map when relying on its own odometry engine, is highly dependent on the environment, hard to parametrize, and is prone to "empty space misconception".

Odometry	KITTI Sequence											time(msec)
	00	01	02	03	04	05	06	07	08	09	10	
F2F	1.4	14.5	4.7	0.4	0.2	0.72	1.8	0.6	5.8	2.2	3.0	61
F2M	1.0	4.7	4.7	0.3	0.2	0.5	0.8	0.5	3.8	2.8	0.8	82
ORB-RTAB	1.0	5.3	4.4	0.2	0.2	0.5	0.6	0.5	3.0	1.5	0.9	175
ORB-SLAM2	1.3	10.4	5.7	0.6	0.2	0.8	0.8	0.5	3.6	3.2	1.0	-

Table 1: Baseline for KITTI Dataset

4 Research Methodology and Procedure

4.1 Baseline Evaluation

The baseline is the accuracy of RTAB-Map and ORB-SLAM2 on KITTI and TUM datasets. Table 1 presents average translational error(ATE) results from multiple papers for KITTI[1]. The time measurement depends on the computer configurations, thus we will measure the time again on our own configurations. The TUM RGB-D dataset was recorded using a handheld Kinect v1 in small office-like environments[14]. For TUM RGB-D dataset, ATE results can be gained from the RTAB-Map paper[1].

4.2 Depth Information Generation from ORB-SLAM2

RTAB-Map team incorporate ORB-SLAM2 by using ORB-SLAM2 as a odometry input, and disabling loop closing and full bundle adjustment of ORB-SLAM2[1]. Local bundle adjustment of ORB-SLAM2 is still working, which make the modified ORB-SLAM2 similar to F2M.

Because ORB-SLAM2 optimize the feature points in the map, it can provide us with more accurate depth information for the feature points[3] in keyframes and more accurate transformation between different camera poses. In our project, we will utilize depth information after optimization, to perform a more accurate odometry estimation.

The RTAB-Map is a Motion-Only Bundle Adjustment, it optimize camera pose by minimizing the reprojection error between matched 3D points in the world coordinate and the observation of camera:

$$\{\mathbf{R}, \mathbf{t}\} = \operatorname{argmin}_{\mathbf{R}, \mathbf{t}} \sum_{i \in \mathcal{K}} \rho(\|x^i - \pi(\mathbf{R}\mathbf{X}^i + \mathbf{t})\|^2)$$

In ORB-SLAM2, the system optimize the depth information of feature points in keyframes by utilizing keypoints observation in all frames[3]:

$$\{\mathbf{X}^i, \mathbf{R}_1, \mathbf{t}_1\} = \operatorname{argmin}_{\mathbf{X}^i, \mathbf{R}_1, \mathbf{t}_1} \sum_{k \in \mathcal{K}_L \cup \mathcal{K}_F} \sum_{i \in \mathcal{K}_k} \rho(\|x^j - \pi(\mathbf{R}_k \mathbf{X}^j + \mathbf{t}_k)\|^2)$$

Comparing to the implementation of Labbé et. al[1], we will utilize optimized keyframe further from ORB-SLAM2, to gain a better accuracy in both trajectory and point cloud map.

4.3 System Validation

In order to have a unified procedure for evaluating the system performance, publicly available datasets with ground truth are utilized to. The datasets were selected so that they can exploit the system boundaries in a wide spectrum of scenarios and conditions. For that purpose, KITTI dataset [15] was selected to cover scenarios in outdoor conditions with random moving objects, and TUM RGB-D dataset [14] was selected to cover more static and indoor environment conditions.

5 Performance Evaluation

The three corner stones of performance in computer science are accuracy, efficiency, and storage requirements. In this work, we focus on addressing *accuracy* aspects of the system which are defined by the relation between the generated trajectory and ground truth data.

The *Absolute Trajectory Error (ATE)* was defined in [5] as a performance metric to stand on the accuracy of a SLAM system vs. a ground truth. The ATE is defined by:

$$ATE(t_i) = \|(x_{t_i}^*, y_{t_i}^*) - (x_{t_i}, y_{t_i})\| \quad (1)$$

where $(x_{t_i}^*, y_{t_i}^*)$ are the ground truth coordinates at t_i , and (x_{t_i}, y_{t_i}) are the coordinates generated by the SLAM algorithm under test at the same epoch. For a more detailed representation of this performance metric, statistical functions are applied to this metric such as Root Mean Square (RMS), mean, median, variance, and standard deviation. Additionally, the maximum and minimum *ATE* are reported to define the system limits while operating.

6 Significance of Proposed Research

As discussed in previous sections, both RTAB-map and ORB-SLAM2 are state-of-the-art algorithms and are able to provide accurate localization and mapping information under a number of conditions. Integration of both algorithms can extend the boundaries of each and can yield a more robust and stable SLAM system with the ability to work in even more challenging conditions. Such improvement can contribute in the all time quest of having an out-of-the-box SLAM system [6] that can be accurate and optimized for long-term operation.

References

- [1] M. Labbé and F. Michaud, “Rtab-map as an open-source lidar and visual simultaneous localization and mapping library for large-scale and long-term online operation,” *Journal of Field Robotics*, vol. 36, no. 2, pp. 416–446, 2019.
- [2] 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.
- [3] 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.
- [4] A. R. Gaspar, A. Nunes, A. Pinto, and A. Matos, “Comparative study of visual odometry and slam techniques,” in *ROBOT 2017: Third Iberian Robotics Conference* (A. Ollero, A. Sanfeliu, L. Montano, N. Lau, and C. Cardeira, eds.), (Cham), pp. 463–474, Springer International Publishing, 2018.
- [5] M. Filipenko and I. Afanasyev, “Comparison of various slam systems for mobile robot in an indoor environment,” in *2018 International Conference on Intelligent Systems (IS)*, pp. 400–407, Sep. 2018.
- [6] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, “Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age,” *IEEE Transactions on robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [7] G. Bresson, Z. Alsayed, L. Yu, and S. Glaser, “Simultaneous localization and mapping: A survey of current trends in autonomous driving,” *IEEE Transactions on Intelligent Vehicles*, vol. 2, no. 3, pp. 194–220, 2017.
- [8] T. Taketomi, H. Uchiyama, and S. Ikeda, “Visual slam algorithms: a survey from 2010 to 2016,” *IPSJ Transactions on Computer Vision and Applications*, vol. 9, no. 1, p. 16, 2017.
- [9] H. Strasdat, J. M. Montiel, and A. J. Davison, “Visual slam: why filter?,” *Image and Vision Computing*, vol. 30, no. 2, pp. 65–77, 2012.
- [10] F. Fraundorfer and D. Scaramuzza, “Visual odometry: Part ii - matching, robustness, and applications,” *IEEE Robotics & Automation Magazine*, vol. 19, no. 2, pp. 78–90, 2012.
- [11] N. Ragot, R. Khemmar, A. Pokala, R. Rossi, and J. Ertaud, “Benchmark of visual slam algorithms: Orb-slam2 vs rtab-map*,” in *2019 Eighth International Conference on Emerging Security Technologies (EST)*, pp. 1–6, July 2019.

- [12] B. M. F. da Silva, R. S. Xavier, T. P. do Nascimento, and L. M. G. Gonsalves, “Experimental evaluation of ros compatible slam algorithms for rgb-d sensors,” in *2017 Latin American Robotics Symposium (LARS) and 2017 Brazilian Symposium on Robotics (SBR)*, pp. 1–6, Nov 2017.
- [13] I. Z. Ibragimov and I. M. Afanasyev, “Comparison of ros-based visual slam methods in homogeneous indoor environment,” in *2017 14th Workshop on Positioning, Navigation and Communications (WPNC)*, pp. 1–6, Oct 2017.
- [14] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, “A benchmark for the evaluation of rgb-d slam systems,” in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 573–580, IEEE, 2012.
- [15] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.