

Towards Semantic KinectFusion

Nicola Fioraio, Gregorio Cerri, and Luigi Di Stefano

Dept. of Computer Science and Engineering
University of Bologna, viale Risorgimento, 2, Bologna (Italy),
{nicola.fioraio, luigi.distefano}@unibo.it,
WWW home page: <http://vision.deis.unibo.it>

Abstract. In this paper we propose an extension to the KinectFusion approach which enables both SLAM-graph optimization, usually required on large looping routes, as well as discovery of semantic information in the form of object detection and localization. Global optimization is achieved by incorporating the notion of keyframe into a KinectFusion-style approach, thus providing the system with the ability to explore large environments and maintain a globally consistent map. Moreover, we integrate into the system our recent object detection approach based on a new Semantic Bundle Adjustment paradigm, thereby achieving joint detection, tracking and mapping. Although our current implementation is not optimized for real-time operation, the principles and ideas set forth in this paper can be considered a relevant contribution towards a Semantic KinectFusion system.

Keywords: KinectFusion, semantic SLAM, semantic bundle adjustment, object detection.

1 Introduction

In the last decade SLAM (Simultaneous Localization and Mapping) has witnessed impressive progresses [1,2,3,4]. Feature-based approaches [1,2] have been deeply explored, so as to attain both accurate maps and real-time operation. 3D reconstruction, though, turns out typically quite sparse and hence often unsuited to robotic tasks such as motion planning and obstacle avoidance. Recently, novel dense approaches using all image pixels to infer a 3D model of the environment from the video sequence have been proposed [5,6]. Depth cameras, such as the Microsoft Kinect, have then brought in new sensing modalities providing 3D measurements at 30Hz, thus paving the way for a brand-new RGBD-SLAM research topic. RGB-D Mapping, proposed by Henry *et al.* [7], was one of the first RGBD-SLAM frameworks whereby both color and depth information are deployed for tracking. The system also includes loop closure detection and a 3D surface reconstruction engine. Later, Newcombe *et al.* [3] introduced KinectFusion, a real-time tracking algorithm capable of unrivaled accuracy. KinectFusion constantly updates a volumetric representation of the scene consisting of a Truncated Signed Distance Function (TSDF) [8] and tracks incoming frames based on the current 3D reconstruction, achieving real-time operation through efficient GPU implementation. As the main limitation of KinectFusion is the bounded mapped scene volume, Whelan *et al.* [4] proposed to shift the TSDF cube according to camera movements throughout

the environment. While the camera moves, the voxels exiting the current volume are triangulated so that a highly detailed 3D reconstruction of large environments can be obtained. However, even though KinectFusion is a low-drift tracking system, mapping large areas is inherently prone to propagation and amplification of pose errors which, afterward, may lead to a poor-quality reconstruction. As highlighted by the authors, possible solutions to address this issue might be reintegration into the moving volume of previously extracted surfaces, as well as handling loop closure to carry out a global optimization.

In this paper we propose first a novel extension to KinectFusion-style system which uses a TSDF volume to build a keyframe representation of the environment. This approach allows for seamless reintegration of already mapped areas and global bundle-adjustment optimization. Moreover, as a second contribution, we enrich the 3D reconstruction with semantic information by means of Semantic Bundle Adjustment [9]: while the camera moves object instances are detected leveraging SLAM and their poses estimated altogether with camera poses. Though other systems allows for semantic SLAM [10,11] or semantic SFM [12], [9] casts the problem as a bundle adjustment style optimization, enabling a fully integrated pipeline. As a result, our system can close seamlessly medium-size loops, delivers accurate 3D reconstructions and peculiarly provides semantic knowledge concerning the mapped environment.

Next section discusses our novel KinectFusion extension, while in Sec. 3 we will show how to integrate the semantic bundle adjustment framework [9] into the proposed KinectFusion extension. Sec. 4 reports both quantitative and qualitative results, in particular by comparing our methods to RGB-D Mapping. Finally, in Sec. 5 some concluding remarks are drawn.

2 Bundle Adjustment by the TSDF

The TSDF representation of a 3D volume consists of a voxel data structure with each element storing the signed distance to the nearest surface. According to Curless *et al.* [8], a depth image is merged into the current volume by back-projection of voxels onto the image plane; then, the resulting pixel gives the position of the surface and the required distance can be computed. Denoted as $[F_{k-1}(\mathbf{p}), W_{k-1}(\mathbf{p})]$, respectively, the signed distance value and a corresponding weight at the 3D location \mathbf{p} at time $k-1$, a new measurement $[F(\mathbf{p}), W(\mathbf{p})]$, coming from the projection of \mathbf{p} onto the image plane is integrated as:

$$F_k(\mathbf{p}) = \frac{W_{k-1}(\mathbf{p})F_{k-1}(\mathbf{p}) + W(\mathbf{p})F(\mathbf{p})}{W_{k-1}(\mathbf{p}) + W(\mathbf{p})} \quad (1)$$

$$W_k(\mathbf{p}) = W_{k-1}(\mathbf{p}) + W(\mathbf{p}) \quad (2)$$

If the signed distance appears to be over a threshold μ from the surface, it must be truncated, as discussed in [8]. In the following, we will refer to the voxels having a non-truncated function value as *surface voxels*.

Recently, the KinectFusion system by Newcombe *et al.* [3] showed how to exploit the TSDF representation not only to achieve detailed surface reconstruction, but also

during the camera tracking process. First, a synthetic depth map is rendered by ray-casting each image pixel according to the known camera intrinsic parameters and following the ray to the intersection with the surface. As described in [8,3], this process is sped up by the TSDF, as the surface is implicitly represented by the zero-crossings of the function. Then, each incoming frame is registered to this virtual depth image through an ICP-like scheme [13,14] and the final aligned frame merged with the surface according to the estimated pose.

The continuous process of TSDF ray-casting and updating makes the system robust and produces low-drift mapping, especially for very loopy trajectories. However, the extent of the mappable volume is bounded by the initial TSDF volume so that it is not possible to reconstruct larger environments. While Whelan *et al.* [4] proposed to simply move the TSDF volume, we think that this strategy is inherently unable to maintain a globally consistent map, since it is not clear how previously mapped locations can be taken into account for global error minimization, *e.g.* how to achieve consistent mapping across large looping routes.

Following a strategy common to many successful SLAM algorithms [7,15], we introduce the notion of *keyframe* as a spatial sampling of camera trajectory and globally optimize keyframe poses for consistent mapping. However, unlike previous proposals, keyframe selection and constraining is performed through the TSDF representation of the environment. To achieve this, we add a new field to each TSDF voxel to keep track of a list of keyframes. Then, when a new frame promoted to keyframe is merged into the volume, it is also added to the keyframe list of all its surface voxels, *i.e.* those voxels to whom a non-truncated distance value has been assigned. Accordingly, on one hand for every point close to the estimated surface we always know which keyframes it was captured in. On the other hand, surface voxels with no associated keyframes clearly represent non-mapped space. Therefore, at the end of the merging phase, we can mark a frame as a keyframe if the ratio between non-mapped surface voxels and total number of surface voxels in the TSDF is above a threshold, *i.e.* the frame describes a volume of the explored environment only partially overlapped with the current map.

Every time a new keyframe is detected and merged, a new optimization step is carried out to ensure consistency over long trajectories. Purposely, a pose graph is created and a suitable cost function defined over this graph structure is minimized by the G2O optimizer [16]. For each unknown keyframe pose we add a vertex to the graph and link it to the other vertex poses with edges representing constraints in the form of non-linear least-squares terms:

$$e_{ab} = \|\mathbf{p}_a - \mathbf{T}_a^{-1}\mathbf{T}_b\mathbf{p}_b\|^2 \quad (3)$$

with $\mathbf{p}_a, \mathbf{p}_b \in \mathbb{R}^3$ representing two corresponding 3D points from keyframe a and b , $\mathbf{T}_a, \mathbf{T}_b$ the unknown poses of those frames. Although data association may be performed densely, for the sake of computational efficiency we overlay a regular grid onto the image and match only such sampled points [17]. Fig. 1 visually summarizes our novel matching scheme; in particular, for each sampled point we:

1. find its corresponding surface voxel by back-projection;
2. get the keyframe list of that voxel and project the 3D point onto all those image planes;
3. perform a local 2D search and select the nearest 3D point under some conditions.

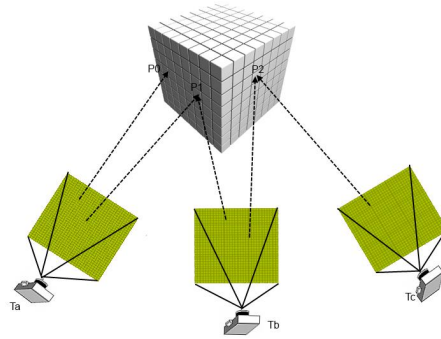


Fig. 1: For each point sampled on keyframes T_a , T_b , T_c we found its matches by projection onto all image planes relative to the keyframe list of the corresponding surface voxel

It is noteworthy to point out that the keyframe list stored in every surface voxel allows for searching only in a well-defined subset of possible keyframes, thus achieving both higher matching quality as well as faster speed. Moreover, the local search in the proximity of the projection ensures to find the nearest 3D point, possibly fulfilling also other requirements, such as e.g. alignment of surface normals [18]. The resulting set of vertex poses and constraints is finally optimized so as to achieve global error minimization.

Once keyframe poses have been jointly optimized, the TSDF no longer represents the mapped environment and it must therefore be recreated to reflect the new estimate of the trajectory. To allow for unbounded camera movements, the new reconstruction is centered at the last keyframe's camera pose. Though definitely more expensive, we claim that such volume shifting strategy is more effective than the proposal in [4] because:

- we do not re-center the active volume only looking at the translation part of the last estimated pose, but quantitatively evaluating the novelty of the brought information;
- we ensure consistent mapping by detecting keyframes to be optimized.

Finally, to reduce the computational effort, the optimization is usually limited to those keyframes having their camera reference frame inside the active volume. This approach was inspired by the Relative Bundle Adjustment framework [19].

3 Semantic KinectFusion

A standard SLAM approach would not deploy any kind of semantic information, such as e.g. the presence of known objects within the explored environment (see Fig. 2a), although this may set forth useful constraints concerning camera poses. Likewise, a major nuisance hindering standard object detection approaches in cluttered scenes deals with having to establish upon objects' presence and pose based on a single possibly unfavorable viewpoint, whereas a SLAM framework would in principle enable continuous incremental discovery from several known viewpoints (see Fig. 2b). In essence, unlike

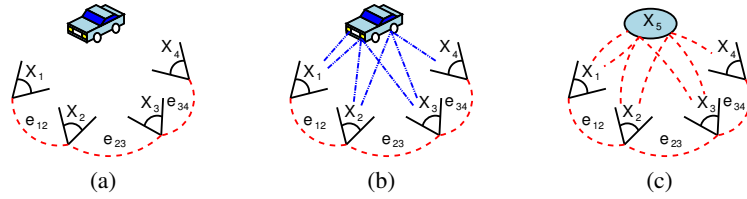


Fig. 2: A standard pose graph (a) ignores any semantic information. Instead, we include into the optimization object matches (b) as graph edges and the object pose as a vertex (c) so as to achieve object detection and improve SLAM.

the state-of-the-art paradigms in both fields, we envision SLAM and object detection as tightly connected and synergistic rather than disjoint processes. Accordingly, we have recently proposed a novel Semantic Bundle Adjustment framework [9], whereby unknown object poses are explicitly included into the graph as pose vertexes constrained to frames by a set of verified hypotheses (see Fig. 2c). Our proposal uniquely attain semantically constrained SLAM together with multi-view object detection and 6DOF localization.

When a frame is promoted to keyframe, 3D features, *e.g.* [20,21], are extracted and matched against those stored into a database of object models. Then, a set of candidate hypotheses on objects' presence and poses is drawn by RANSAC-based 6DOF pose estimation [22] and a *validation graph* is created to verify the consistency of each of such hypotheses with respect to the current 3D reconstruction. As detailed in [9], each validation graph is populated with all the vertex poses featuring object matches, their edges, a new vertex pose for the unknown object location and an edge for each matching feature:

$$e_{fk}^n = s_{fk}^n \left\| \mathbf{p}_f - \mathbf{T}_f^{-1} \mathbf{T}_k \mathbf{p}_k^n \right\|^2 \quad (4)$$

\mathbf{p}_k^n being the n^{th} 3D point feature on object k , \mathbf{p}_f the matching 3D point feature on frame f , s_{fk}^n the matching score, \mathbf{T}_f and \mathbf{T}_k the unknown frame and object poses. Moreover, virtual edges are created when different camera frames match the same object feature [9]. Let $\mathbf{p}_{f_0} \in f_0$ and $\mathbf{p}_{f_1} \in f_1$ be two matches for object feature \mathbf{p}_k^n , then, under certain assumptions [9], we can add to the graph the virtual constraint:

$$e_{f_0 f_1}^n = s_{f_0 k}^n s_{f_1 k}^n \left\| \mathbf{p}_{f_0} - \mathbf{T}_{f_0}^{-1} \mathbf{T}_{f_1} \mathbf{p}_{f_0} \right\|^2 \quad (5)$$

The validation graph is optimized and a cleaning procedure is performed to detect wrong hypotheses (see [9] for more details). Finally, if the detection is confirmed, all the constraints are included into the global graph (see Sec. 2) and a global semantic optimization is performed.

Fig. 3 summarizes our overall proposal. Each incoming frame is fed to the tracking module (left box), which merges the new measurements with the current TSDF volume and checks the size of the non-mapped area to detect keyframes. If a new keyframe is spawned, 3D features are extracted and matched against models' features (right box),

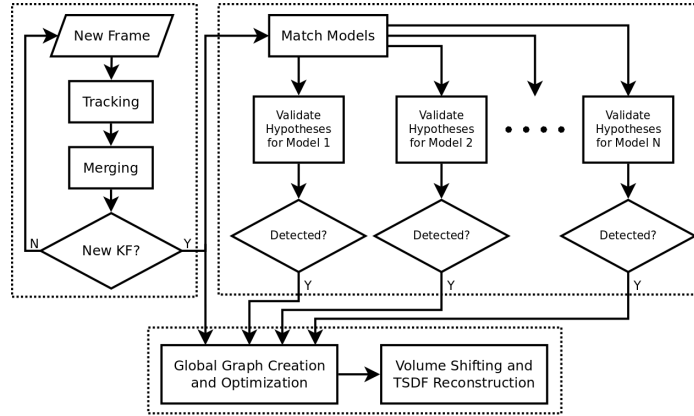


Fig. 3: A simplified view of our system.

thus yielding hypotheses which have to be verified by the semantic bundle adjustment process. Finally, detected objects and keyframes are jointly optimized (bottom box) and the TSDF volume is updated.

4 Results

The proposed approach has been extensively tested, both quantitatively as well as qualitatively. As for quantitative experiments, to the best of our knowledge available RGB-D datasets, proposed either for SLAM or 3D object detection, cannot be used to assess joint SLAM *and* 3D object detection, as the ground truth concerns camera trajectory only (in SLAM datasets) or, alternatively, object identities and poses only (in object detection datasets). Therefore, as proposed in [9], we rely here on a semi-synthetic setup: 3D object models scanned by a Kinect are rendered by ray-casting into SLAM datasets according to ground-truth camera poses. Although object detection turns out easier than in a real environment, we found such semi-synthetic datasets to deliver valuable insights and feedbacks concerning behavior and accuracy of the proposed method.

Hence, we took the video sequences from the RGB-D SLAM dataset [15], using the provided ground truth trajectory for precise model placement and ray-casting. Then, we also captured sequences by a Kinect camera featuring known objects in real environments, so as to provide also qualitative results on truly real data. As for the feature matching step underlying object detection, in all experiments we used the SIFT3D key-point detector [23,24] and the Color-SHOT feature descriptor [21], using for both algorithms the implementations publicly available in the open-source Point Cloud Library [24].

Tab. 1 and 2 report the results obtained, respectively, with our implementation of RGB-D Mapping and with the extension to KinectFusion described in Sec. 2. As both such systems do not incorporate semantic information into the SLAM process, for their quantitative evaluation we can use the original RGB-D SLAM sequences [15]. As vouched by the Tables, we usually perform better, especially when considering the

Table 1: Quantitative results with RGB-D Mapping (translation error wrt ground-truth camera poses). *Kfs* denotes the error at keyframes only, *All Frs* the error for all estimated camera poses.

Sequence	RGB-D Mapping					
	<i>Mean Err. (m)</i>		<i>Max Err. (m)</i>		<i>RMSE (m)</i>	
	Kfs	All Frs	Kfs	All Frs	Kfs	All Frs
FR1 Floor	0.056	0.093	0.1	0.244	0.06	0.104
FR1 360	0.089	0.197	0.211	0.4	0.099	0.208
FR1 Desk	0.069	0.074	0.163	0.181	0.076	0.081
FR1 Desk2	0.111	0.085	0.219	0.262	0.121	0.097

Table 2: Quantitative results with our extension to KinectFusion (translation error wrt ground-truth camera poses). *Kfs* denotes the error at keyframes only, *All Frs* the error for all estimated camera poses.

Sequence	Our Extension to KinectFusion					
	<i>Mean Err. (m)</i>		<i>Max Err. (m)</i>		<i>RMSE (m)</i>	
	Kfs	All Frs	Kfs	All Frs	Kfs	All Frs
FR1 Floor	0.056	0.054	0.093	0.13	0.062	0.058
FR1 360	0.083	0.183	0.159	0.324	0.091	0.192
FR1 Desk	0.04	0.055	0.131	0.167	0.047	0.058
FR1 Desk2	0.096	0.083	0.135	0.210	0.098	0.09

Table 3: Quantitative results with our Semantic KinectFusion (translation error wrt ground-truth camera poses). *Kfs* denotes the error at keyframes only, *All Frs* the error for all estimated camera poses.

Sequence	Our Semantic KinectFusion					
	<i>Mean Err. (m)</i>		<i>Max Err. (m)</i>		<i>RMSE (m)</i>	
	Kfs	All Frs	Kfs	All Frs	Kfs	All Frs
FR1 Floor	0.054	0.05	0.092	0.106	0.059	0.055
FR1 360	0.066	0.144	0.167	0.249	0.073	0.152
FR1 Desk	0.044	0.069	0.092	0.169	0.048	0.076
FR1 Desk2	0.042	0.085	0.088	0.132	0.046	0.088

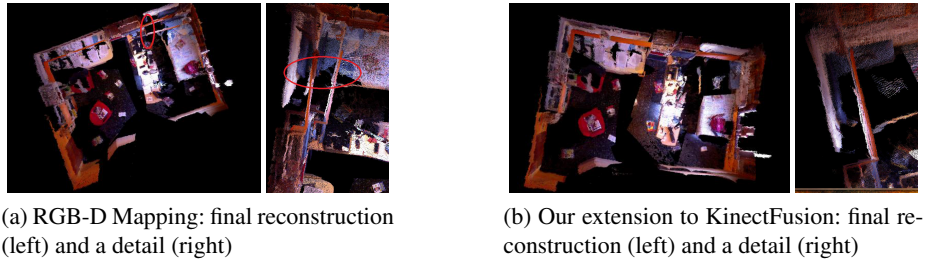


Fig. 4: We moved the camera on a circular trajectory and repeatedly acquired the same loop. Our proposal always tracks against the map, RGB-D mapping eventually drifts away.

whole trajectory (*All Frs* column). Indeed, when the sensor comes back to an already mapped place, previous keyframes correctly enters the active TSDF volume thanks to the small drift and we can start tracking with respect to the known map. Instead, RGB-D Mapping looks for loop closure every time a new keyframe is detected, with no limit to their proliferation. Such effect is particularly evident in Fig. 4, which comes from the dataset used for qualitative evaluation (*i.e.* our own dataset without ground-truth information). Here camera motion follows a circular trajectory performing several loops around the room: at the end of the first loop, our technique hangs on to the known map, reducing the localization error and thus not affecting the reconstruction due to no new keyframe being spawned.

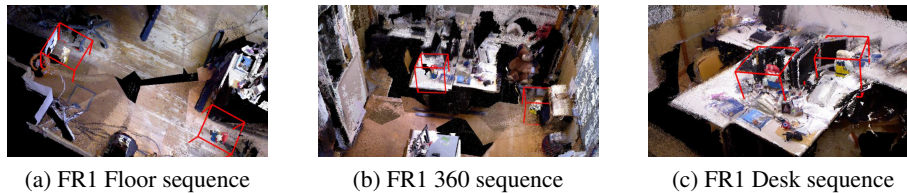


Fig. 5: Semi-synthetic RGB-D SLAM dataset: the semantic framework achieves high quality reconstruction and accurate object detection and localization.

As for our proposed Semantic KinectFusion approach, Tab. 3 shows no reduction of accuracy, and often some improvements due to semantic information effectively constraining camera poses across many views. Moreover, beside a high quality reconstruction, the semantic framework can precisely localize the objects looked for, as shown in Fig. 5. Finally, as object detection is easier with the semi-synthetic setup than in real settings, we show also qualitative experiments on real data acquired by a Kinect camera in indoor environments. Fig. 6 reports the results obtained adding physically an object belonging to the model database into the same room as in Fig. 4. From one hand,

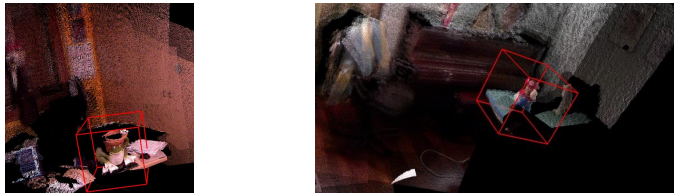


Fig. 6: We moved the camera on a circular trajectory and repeatedly acquired the same loop. TSDF-enabled optimization yields very low drift tracking, while semantic bundle adjustment detects and localizes objects (see the red boxes).

our novel TSDF-enabled optimization (Sec. 2) achieves high mapping accuracy and low drift even on multiple loops; on the other hand, the Semantic Bundle Adjustment framework (Sec. 3) detects and localizes the object in the scene, thereby providing also semantic constraints which help improving the overall global alignment.

5 Concluding Remarks

A novel TSDF-based technique for bundle adjustment style optimization has been presented. Constraints across views are obtained by means of the TSDF representation of the environment and keyframes, alike, are selected with respect to such a representation. Then, we integrated a state-of-the-art Semantic Bundle Adjustment framework into the system, thus achieving effective joint detection, tracking and mapping, as vouched by both quantitative and qualitative experiments.

Future works will definitely concern real-time implementation on a modern GPU architecture, so as to extend the successful KinectFusion framework with global bundle adjustment and seamless object detection and localization. We also plan to investigate on whether and how effective features for object detection may be extracted directly from the TSDF reconstruction, which holds the potential to provide a more comprehensive and robust description than a single camera frame, as well as on how the TSDF-enabled BA could be extended to objects too. Finally, in this work we chose keyframes among camera frames; however, an interesting alternative we wish to explore deals with rycasting of the scene to get dense, low noise, virtual keyframes.

References

1. A. J. Davison, “Real-time simultaneous localisation and mapping with a single camera,” in *Computer Vision (ICCV), IEEE Int’l Conf. on*, (Washington, DC, USA), p. 1403, 2003.
2. G. Klein and D. Murray, “Parallel tracking and mapping for small ar workspaces,” in *Mixed and Augmented Reality (ISMAR), IEEE and ACM Int’l Symp on*, pp. 225–234, nov. 2007.
3. R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon, “Kinectfusion: Real-time dense surface mapping and tracking,” in *ISMAR*, (Washington, DC, USA), pp. 127–136, 2011.
4. T. Whelan, J. McDonald, M. Kaess, M. Fallon, H. Johannsson, and J. Leonard, “Kintinuuous: Spatially extended KinectFusion,” in *RSS Workshop on RGB-D: Advanced Reasoning with Depth Cameras*, (Sydney, Australia), Jul 2012.

5. J. Stühmer, S. Gumhold, and D. Cremers, “Real-time dense geometry from a handheld camera,” in *Proceedings of the 32nd DAGM conference on Pattern recognition*, (Berlin, Heidelberg), pp. 11–20, Springer-Verlag, 2010.
6. R. Newcombe, S. Lovegrove, and A. Davison, “Dtam: Dense tracking and mapping in real-time,” in *Computer Vision (ICCV), IEEE Int’l Conf. on*, pp. 2320–2327, nov. 2011.
7. P. Henry, M. Krainin, E. Herbst, X. Ren, and D. Fox, “Rgb-d mapping: Using depth cameras for dense 3d modeling of indoor environments,” in *Proc. of Int’l Symp on Experimental Robotics (ISER)*, 2010.
8. B. Curless and M. Levoy, “A volumetric method for building complex models from range images,” in *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques, SIGGRAPH ’96*, (New York, NY, USA), pp. 303–312, ACM, 1996.
9. N. Fioraio and L. Di Stefano, “Joint detection, tracking and mapping by semantic bundle adjustment,” in *Computer Vision and Pattern Recognition (CVPR), IEEE Conf. on*, (Portland, OR, USA), 2013.
10. J. Civera, D. Gálvez-López, L. Riazuelo, J. D. Tardós, and J. M. M. Montiel, “Towards semantic SLAM using a monocular camera,” in *Proc. of the Int’l Conf. on Intelligent Robot Systems (IROS)*, pp. 1277–1284, 2011.
11. S. Ekvall, P. Jensfelt, and D. Kragic, “Integrating active mobile robot object recognition and slam in natural environments,” in *Intelligent Robots and Systems, IEEE/RSJ Int’l Conf. on*, oct. 2006.
12. S. Y. Bao and S. Savarese, “Semantic structure from motion,” in *CVPR*, 2011.
13. Y. Chen and G. Medioni, “Object modelling by registration of multiple range images,” in *Proc. of the IEEE Int’l Conf. on Robotics and Automation*, vol. 3, pp. 2724–2729, April 1991.
14. P. J. Besl and H. D. McKay, “A method for registration of 3-d shapes,” *PAMI*, vol. 14, no. 2, pp. 239–256, 1992.
15. F. Endres, J. Hess, N. Engelhard, J. Sturm, D. Cremers, and W. Burgard, “An evaluation of the RGB-D SLAM system,” in *Robotics and Automation (ICRA), IEEE Int’l Conf. on*, (St. Paul, MA, USA), May 2012.
16. R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, “g2o: A general framework for graph optimization,” in *ICRA*, (Shanghai, China), may 2011.
17. N. Fioraio and K. Konolige, “Realtime visual and point cloud slam,” *Proc. of the RGB-D Workshop on Advanced Reasoning with Depth Cameras at Robotics: Science and Systems Conf.(RSS)*, vol. 27, 2011.
18. S. Rusinkiewicz and M. Levoy, “Efficient variants of the ICP algorithm,” in *Proceedings Third International Conference on 3-D Digital Imaging and Modeling*, pp. 145–152, IEEE Comput. Soc, 2001.
19. G. Sibley, C. Mei, I. Reid, and P. Newman, “Adaptive relative bundle adjustment,” in *Robotics Science and Systems (RSS)*, (Seattle, USA), june 2009.
20. A. Johnson, *Spin-Images: A Representation for 3-D Surface Matching*. PhD thesis, Robotics Institute, Carnegie Mellon University, August 1997.
21. F. Tombari, S. Salti, and L. Di Stefano, “A combined texture-shape descriptor for enhanced 3D feature matching,” in *18th IEEE Int’l Conf. on Image Processing (ICIP)*, (Brussels, Belgium), pp. 809–812, September, 11-14 2011.
22. K. S. Arun, T. S. Huang, and S. D. Blostein, “Least-squares fitting of two 3-d point sets,” *Pattern Analysis and Machine Intelligence (PAMI), IEEE Trans. on*, vol. 9, pp. 698–700, sept 1987.
23. D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *IJCV*, vol. 60, pp. 91–119, January, 5 2004.
24. R. B. Rusu and S. Cousins, “3D is here: Point cloud library (PCL),” in *IEEE Int’l Conf. on Robotics and Automation (ICRA)*, (Shanghai, China), May 9-13 2011.