

RGB-D SLAM with Structural Regularities

Yanyan Li¹, Raza Yunus¹, Nikolas Brasch¹, Nassir Navab^{1,2} and Federico Tombari^{1,3}

Abstract—This work proposes a RGB-D SLAM system specifically designed for structured environments and aimed at improved tracking and mapping accuracy by relying on geometric features that are extracted from the surrounding. Structured environments offer, in addition to points, also an abundance of geometrical features such as lines and planes, which we exploit to design both the tracking and mapping components of our SLAM system. For the tracking part, we explore geometric relationships between these features based on the assumption of a Manhattan World (MW). We propose a decoupling-refinement method based on points, lines, and planes, as well as the use of Manhattan relationships in an additional pose refinement module. For the mapping part, different levels of maps from sparse to dense are reconstructed at a low computational cost. We propose an instance-wise meshing strategy to build a dense map by meshing plane instances independently. The overall performance in terms of pose estimation and reconstruction is evaluated on public benchmarks and shows improved performance compared to state-of-the-art methods. The code is released at <https://github.com/yanyan-li/PlanarSLAM>.

I. INTRODUCTION

Visual Simultaneous Localization and Mapping (SLAM) algorithms are used to estimate the 6D camera pose while reconstructing the surrounding unknown environment. They have shown to be useful in a wide range of applications, such as autonomous robots, self-driving cars and augmented/virtual reality, where camera pose estimation enables cars, robots and mobile devices to localize themselves, while the dense map provides a representation of the environment, *e.g.* for robot-environment or human-environment interaction.

Many SLAM applications have to deal with structured scenes, *i.e.* man-made environments that are usually characterized by low-textured surfaces - a typical example is an indoor scene, or an outdoor parking place. This induces a lack of visual features, that visual SLAM systems typically leverage to improve camera pose estimation and/or 3D reconstruction, *e.g.* by carrying out loop closure and bundle adjustment to reduce drift. In order to deal with structured scenes, specific SLAM methods based on points and line segments, like S-SLAM [1], Stereo-PLSLAM [2], PLVO [3], Mono-PLSLAM [4] and Probabilistic-VO [5] have been proposed, extending the working environment to scenes where more lines than points can be detected. SP-SLAM [6] merges plane features into ORB-SLAM2 [7], achieving robust results in low-textured scenes.

¹Technical University of Munich, Germany; {yanyan.li, raza.yunus, nikolas.brasch, nassir.navab, federico.tombari}@tum.de; ²Johns Hopkins University, USA; ³Google Inc.

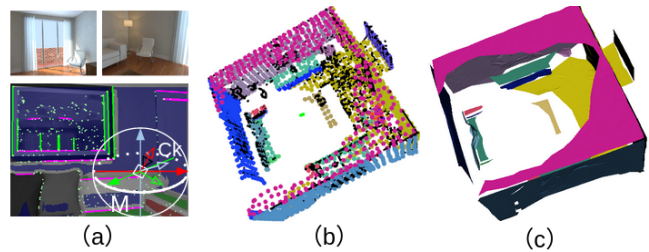


Fig. 1. RGB-D SLAM system. (a) Examples of a typical structured scene, and 2D features and orthogonal lines and planes segmentation. (b) Point cloud including points, lines and planes. (c) Real-time mesh on a CPU.

For the reconstruction, there are sparse, semi-dense and dense methods. Compared to the first two classes, which only provide incomplete maps, dense reconstruction is required to provide sufficient information for applications such as robot-environment interaction and 3D scene understanding. Many algorithms have been proposed to reconstruct indoor scenes via RGB-D sensors. KinectFusion [8] is a pioneering work relying on the truncated signed distance field (TSDF) representation of the map. In order to reconstruct large scale scenarios, surfel-based methods, like ElasticFusion [9], were proposed. Instead of reconstructing each pixel, Wang et al. [10] extracts superpixels from RGB images and depth maps, which is more efficient but still has redundant information especially in indoor scenarios where large planes can be commonly found.

In this paper, we build on our monocular Structure-SLAM [1] and propose a robust RGB-D SLAM system specifically designed to deal with structured environments, which improves tracking and mapping at the same time. Figure 1 illustrates the components of such structured scenes, which contains points, lines and plane segments. Following the decoupling strategy of Structure-SLAM, we estimate a drift-free rotation matrix first, and then the 3-DoF translation. The initial rotation and translation are refined via a map-to-frame strategy. Different to [1], [11], [12], plane features are merged into our Manhattan-based framework, which is used to estimation the initial translation vector and retain Manhattan relationships as constrains in the refinement module. Furthermore, an efficient meshing module is proposed that reconstructs the scene structure based on the obtained planar regions in the sparse map. In summary our contributions are:

- Based on the concept of MW-based decoupled pose estimation, we improve the translation estimation by combining point and line features with planes and an additional pose refinement step with Manhattan relationships.

- We propose a planar instance-wise mesh based reconstruction method generating a compact representation of the environment from a sparse point cloud.
- A general framework for real-time RGB-D SLAM where these components are used to localize and map under structured environments with high accuracy.

We evaluate the performance of our approach in terms of both camera pose estimation and reconstruction on public benchmarks, showing improved performance compared to state-of-the-art methods.

II. RELATED WORK

In the following we review the literature related to RGB-D based SLAM systems as well as methods leveraging structural regularity as the MW assumption.

a) RGB-D SLAM.: In [13], [14] it was proposed to use planes over point features whenever possible, as the averaging over multiple depth measurements reduces the noise significantly. In Dense Planar SLAM [15] surfels belonging to the same planar areas are smoothed by fitting a plane to them and back-projecting the surfels onto the plane. Le *et al.* [16] rely on a scene layout consisting of a ground plane and several walls, and use dynamic programming to infer a sequentially consistent assignment of pixels to planes. In Probabilistic-VO [5], the uncertainties of points, lines and planes are modelled explicitly and used during pose estimation, where points, lines and planes are represented in a uniform framework in [17]. A direct SLAM system combining photo-metric and geometric terms is proposed in DVO-SLAM[18] and extended in CPA-SLAM [19] with global planes, where depth measurements are assigned to the global planes with weights.

b) Dense Reconstruction.: While the aforementioned methods have the goal to estimate precise poses and therefore only maintain a map with the most reliable information, several works have been proposed with the goal to create a complete dense reconstruction of the environment. KinetFusion [8] and ElasticFusion [9] explore dense reconstruction for RGB-D sensors. The first method fuses all depth data into a volumetric dense representation, which is used to track the camera pose using ICP. The size of the map is usually limited in volumetric methods due to memory constraints. Different from KinectFusion, ElasticFusion is a map-centric system that reconstructs surfel-based maps of the environment. In order to decrease the number of surfels in the map, superpixel-based surfels are proposed by [10], which reduces the number of surfels compared with ElasticFusion. Recently BAD-SLAM [20] proposed a direct bundle-adjustment approach for RGB-D SLAM. In [21] a textured mesh is extracted from a dense surfel cloud. A direct mesh based reconstruction approach for RGB-D sensors was proposed in [22].

c) Structural Regularity.: A line of works exploits additional constraints and regularities in the world, to improve the reconstruction performance. In [23] and [24] the authors showed that the rotation estimation error is the main reason for long-term drift.

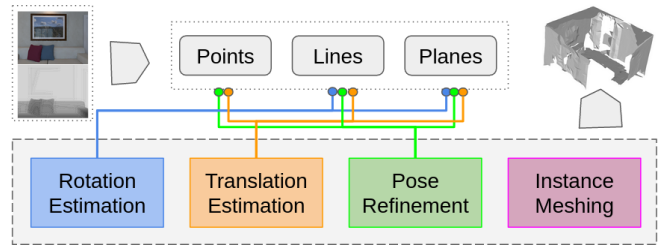


Fig. 2. Overview of the proposed framework. Point, line and plane features are extracted from the RGB-D frame. Rotation and translation are estimated in a decoupled fashion first and refined afterwards. The planar segments are used to create a mesh-based reconstruction of the environment.

A branch-and-bound framework for Manhattan Frame estimation is proposed in [25]. In MVO [24] a method using mean shift on the unit sphere is used to find the transformation between the MW and the current frame. When only planes are used for the rotation estimation as in OPVO [26] at least 2 orthogonal planes must be detected in each frame, the addition of vanishing points extracted from lines can be used alternatively, as done in LPVO [11]. The methods mentioned above use point features to estimate the translation. Structure-SLAM [1] is a monocular system that predicts normals via a convolutional neural network leverages normals with points and lines in a decoupling strategy. Since predicted normals are not as accurate as those computed from a depth map, the system provide a refinement/fallback module based on points and lines. Compared with Structure-SLAM, optimized vanishing points of lines and plane features are used for rotation and translation in this work. Then, the fallback part is removed and the refinement part incorporates geometric relationship of planes. Instead of the sparse point-line map, a dense mesh as output is more useful for robotics applications. L-SLAM [12] is also based on the MW assumption, which obtains translation, rotation and pixels of potential planar regions from LPVO. Then it refines 3D translational and 1-D plane offsets with a linear Kalman Filter. However, we use a more robust front-end for initial translation estimation. Furthermore, the 6D pose refinement step is used to optimize rotation and translation simultaneously and allows an offset to the initial rotation from MW, which is more robust to non-MW (curved surfaces and few planar regions) compared with L-SLAM and LPVO (see Figure 6).

III. PROPOSED FRAMEWORK

Given a sequence of RGB-D frames from a structured environment, the goal of our method is to reconstruct the 3D scene while simultaneously estimating the 6D camera pose at each frame. Section IV provides an overview of the proposed tracking pipeline, which decouples rotation and translation, while section V describes different types of mapping presentations generated by the system. We now describe the system's underlying features and structural components.

A. Extended feature set

In our method we use ORB features [27], which are fast to extract and match. In low-textured environments, it is

hard to extract sufficient points for robust pose estimation, therefore we extend the feature set with lines, which are extracted using the LSD [28] approach, as described by LBD [29]. Furthermore, it is common to find non-textured planar regions in indoor environments, where plane instances extracted from the depth maps are valuable cues to extend points and lines. Planes are detected using the connected component analysis method [30]. They are represented by the Hessian normal form $\pi = (\hat{n}, d)$, where $\hat{n} = (n_x, n_y, n_z)$ is the normal of the plane, representing its orientation and d is the distance from the camera origin to the plane.

a) *Points and lines*: After the extraction of 2D point features $x_j = (u_j, v_j)$ and line segments $l_j = (x_{j,start}, x_{j,end})$ in frame F_i , we can back-project points and lines using the camera intrinsic parameters and the depth map to obtain 3D points X_j and 3D lines L_j . The depth map is not always correct, especially at depth discontinuities *e.g.* object boundaries. Therefore a robust fitting method for 3D lines is needed. First, we count the number of pixels with non-zero depth values intersected by the detected line segment. If the number exceeds a certain threshold, the 3D line L_j will be estimated via RANSAC to remove potential outliers.

b) *Normals and planes*: Smooth normals are computed by averaging the tangential vectors from the depth image inside a patch of 10×10 pixels using integral images. After plane detection, we use the strategy of [6] to associate the observed planes with those present in the map. To match an observed plane with one from the map, we first check the angle between their normals. If it is below the threshold θ_n , we check the point-to-plane distance between them. The plane which has the minimum distance to the observed plane, and also lies below the distance threshold θ_P , is matched to the observed plane. In the experiments, θ_n and θ_P are set at 10 degrees and 0.1 m respectively. Furthermore, we also keep parallel and perpendicular relationships [6] between the map planes to leverage additional constraints during the tracking process. These are determined by the angle between the plane normals. Since they only provide constraints for orientation, we do not consider their distance.

B. Decoupling pose estimation and refinement.

To reduce error propagation between frames, we build on our monocular architecture [1] that computes rotational motion based on the MW assumption. Then the corresponding translational motion is estimated by features, with the fixed rotation computed from last step. In the front-end of this work, we use optimized lines for rotation estimation and planes for translation estimation.

Differently to Structure-SLAM [1] that uses a point-line local map to optimize translation and rotation together, we leverage planes in the local map and also make use of the geometric relationship (parallel and perpendicular) of those planes as constraints, which improves the accuracy of the system as it will be shown in Figure 4 and Table I.

IV. TRACKING

Differently from traditional pose estimation methods, we decouple the 6D camera pose into rotation and translation.

Based on the MW assumption, we obtain the rotational motion $R_{c_i m}$ between the MW and camera c_i . In this way, the rotation estimation will not be affected by the pose of the last frame or last keyframe, which reduces drift effectively. Afterwards, point, line and plane features as well as the initial rotation matrix are used for translation estimation, which consists of just 3 Degrees-of-Freedom (DoFs).

A. Rotation estimation

Instead of tracking the camera from frame-to-frame directly, the drift-free rotation estimation method estimates the rotation R_{cm} between each frame and the Manhattan coordinate frame, by modeling the indoor environments as a MW, thus reducing the drift generated from frame-to-frame tracking. As shown in Figure 1, Manhattan coordinate frames can be aligned to the starting frame of the camera via $R_{k+1, m}$. Generally, the coordinate of the first frame is regarded as the world frame, *i.e.* $R_{1, m} = R_{m, w}^T$. So we can obtain pose in the world coordinate by using,

$$R_{k+1, w} = R_{k+1, m} R_{m, w} \quad (1)$$

Here $R_{m, w}$ represents the relation from the world to MW, which is obtained by the MW initialization step and $R_{k+1, m}$ is the relation from MW to the $(k+1)^{th}$ frame. These two matrices are computed via a sphere mean-shift method [24], where the normals and normalized vanishing directions are projected onto the tangent planes of the current rotation estimate. Then a mean shift step is performed on the tangent planes, which generates new centers and back-projects them to the sphere as new estimates. We refer the reader to [24] and [26] for more details on the sphere mean-shift method. To handle difficult scenes where only one or no plane at all is detected, we feed the unit sphere with both vanishing directions of the refined 3D lines and surface normals of planes, which is a more robust approach than [26], [1] under these challenging conditions.

B. Translation estimation

After estimating the rotation, points, lines and planes are used to estimate the translation. We re-project 3D points from the last frame into the current one and define the error function, based on the re-projection error, as follows,

$$e_{k, j}^p = p_k - \Pi(R_{k, j} P_j + t_{k, j}) \quad (2)$$

where $\Pi(\cdot)$ is the projection function. Since the rotation matrix $R_{k, j}$ has been obtained in the last step, we fix the rotation and only estimate the translation using the Jacobian matrix corresponding to (2).

As for lines, we obtain the normalized line function from the 2D endpoints p_{start} and p_{end} as follows

$$l = [p_{start} \times p_{end}] / [\|p_{start}\| \|p_{end}\|] = (a, b, c). \quad (3)$$

Then, we formulate the error function based on the point-to-line distance [4] between l and the projected 3D endpoints P_{start} and P_{end} from the matched 3D line in the keyframe. For each endpoint P_x , the error function can be noted as,

$$e_{k, P_x}^l = l \Pi(R_{k, j} P_x + t_{k, j}). \quad (4)$$

To get a minimal parameterization of a plane π for optimization, we represent it as $q(\pi) = (\phi, \psi, d)$ where ϕ and ψ are the azimuth and elevation angles of the normal and d is the distance from the Hessian form

$$q(\pi) = (\phi = \arctan(\frac{n_y}{n_x}), \psi = \arcsin(n_z), d). \quad (5)$$

So, the error function between the observed plane π_k in the frame and corresponding map plane π_x is

$$e_{k,\pi_x}^\pi = q(\pi_k) - q(T_{cw}^{-T} \pi_x) \quad (6)$$

where T_{cw}^{-T} is the transformation from world to camera coordinates. Assuming that the observations follow a Gaussian distribution, the final non-linear least squares cost function t^* can be written as in (7), where $\Lambda_{p_{k,j}}$, $\Lambda_{p_{k,P_x}}$ and Λ_{k,π_x} are the inverse covariance matrices of points, lines and planes, and ρ_p , ρ_l and ρ_π are robust Huber cost functions, respectively.

$$t^* = \underset{j}{\operatorname{argmin}} \sum_j^M \rho_p \left(e_{k,j}^p{}^T \Lambda_{p_{k,j}} e_{k,j}^p \right) + \rho_l \left(e_{k,P_x}^l{}^T \Lambda_{p_{k,P_x}} e_{k,P_x}^l \right) + \rho_\pi \left(e_{k,\pi_x}^\pi{}^T \Lambda_{k,\pi_x} e_{k,\pi_x}^\pi \right) \quad (7)$$

Here, a solution is determined using the Levenberg-Marquardt algorithm.

C. Pose refinement

The last two steps assume that the scene is a good Manhattan model, nevertheless several general indoor environments are not strictly adhering to the MW assumption, leading to degradation in accuracy. So, after obtaining the initial pose via the decoupled rotation and translation strategy, the refinement module [1] fine-tunes the pose to compensate for deviations from the MW or unstable initial estimates. In the refinement step, to reduce the drift from frame-to-frame pose estimation, the local map constructed by previous keyframes is used to optimize the pose based on a map-to-frame strategy [7].

Similar to [6], [7], [31], we also use keyframes to build a local map, although our map has point, line and plane landmarks, which are projected into the current frame to search for matches. Furthermore, we explore the relationship between planes in the local map and planes detected in the current frame. The parallel and perpendicular constraints between those planes are described as (8),

$$\begin{cases} e_{k,n_x}^{\pi_{\parallel}} = ||q_n(n_k) - q_n(R_{cw} n_x)|| \\ e_{k,n_x}^{\pi_{\perp}} = ||q_n(R_{\perp} n_k) - q_n(R_{cw} n_x)|| \end{cases} \quad (8)$$

where $q_n(\pi) = (\phi, \psi)$ and R_{cw} is the transformation from world to camera coordinates. For perpendicular planes, their plane normal is rotated by 90 degrees (R_{\perp}) to construct the error function. These two error functions are merged to (7) to build a joint optimization function in the refinement module.

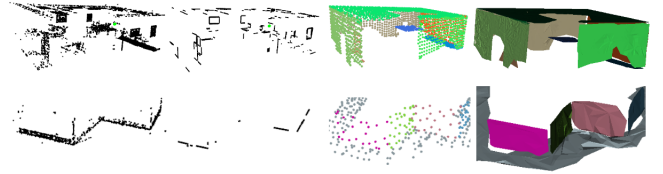


Fig. 3. Different levels of maps provided by the system. Top row: office room of the ICL-NUIM; bottom row: structure-nontexture-near of TUM RGB-D;

V. MAPPING

This section describes the keyframe-based 3D mapping strategy used in our SLAM framework. Keyframes and 3D features build up a co-visibility graph, where nodes and edges are updated whenever a new keyframe and new features are available.

A. Sparse Mapping

As shown in Figure 3, the sparse map module is reconstructed by point-line-plane features extracted from keyframes. The first frame is set as the first keyframe and the global map is initialized by the landmarks thereby detected. When new points, lines and planes are detected in a new keyframe, which are not in the global map, they will be saved to a local map first. Then we check the quality of the landmarks in the local map, and then push reliable landmarks into a global map after culling bad ones. Different to the matching methods for points and lines, for each detected plane in a new keyframe, we first check whether it is associated with a plane in the map using the strategy described in section III. If we find an association, we add the 3D points of the new plane to the associated plane in the global map and filter out redundancies using a voxel grid to get a compact point cloud again. If the incoming plane is not associated to any plane in the global map, we add it to the map as a new plane.

B. Planar instance-wise meshing

The sparse map obtained in the previous section is still not adequate for applications involving robot-environment interactions, but it provides information about planar and non-planar instances. Therefore, we construct a denser map using an instance-wise meshing strategy. Indoor scenes can be divided into planar and non-planar regions. Planar areas like floors, walls and ceiling have often a large extent, however a dense pixel-wise information does not add to the quality and is highly redundant. So instead of using surfel or TSDF, we regard plane regions as instances that include a small and fixed number of elements independently of their size.

In particular, we input plane instances to the meshing module, which meshes them independently. First, the points belonging to a plane are organized as a kd-tree data-structure. Different to unstructured inputs, our method needs less time for searching several nearest neighbors. Then, we use Greedy Surface Triangulation (GST) [33] to build an instance-wise mesh, which is designed to deal with planar surfaces. Note

Sequence	Ours	Ours/-wo	ORB [7]	PS-SLAM [6]	LPVO [12]	L-SLAM [11]	DVO [18]	InfiniTAM [32]
lr-kt0	0.006	0.025	0.025	0.016	0.015	0.012	0.108	×
lr-kt1	0.015	0.036	0.008	0.018	0.039	0.027	0.059	0.006
lr-kt2	0.020	0.053	0.023	0.017	0.034	0.053	0.375	0.013
lr-kt3	0.012	0.059	0.021	0.025	0.102	0.143	0.433	×
of-kt0	0.041	0.068	0.037	0.032	0.061	0.020	0.244	0.042
of-kt1	0.020	0.028	0.029	0.019	0.052	0.015	0.178	0.025
of-kt2	0.011	0.060	0.039	0.026	0.039	0.026	0.099	×
of-kt3	0.014	0.012	0.065	0.012	0.030	0.011	0.079	0.010
snot-far	0.022	0.026	×	0.020	0.075	0.141	0.213	0.037
snot-near	0.025	×	×	0.013	0.080	0.066	0.076	0.022
cabinet	0.035	0.057	0.075	0.067	0.520	0.291	0.690	0.035
large-cabinet	0.071	0.813	0.124	0.079	0.279	0.140	0.979	0.512

TABLE I

COMPARISON OF TRANSLATION RMSE (M) FOR ICL-NUIM AND TUM-RGB-D SEQUENCES. × MEANS THE METHOD FAILS IN THE TRACKING PROCESS. -WO MEANS ONLY USING DECOUPLED TRACKING WITHOUT THE REFINEMENT STEP.

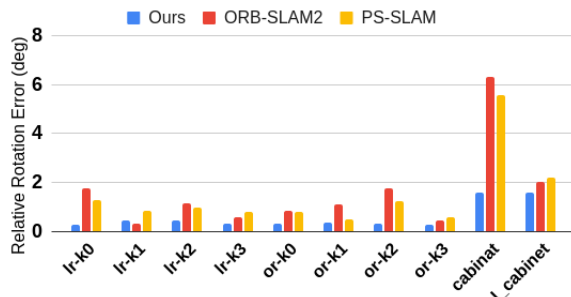


Fig. 4. Comparison of relative pose error (RPE) for rotation on the ICL-NUIM and TUM RGB-D sequences.

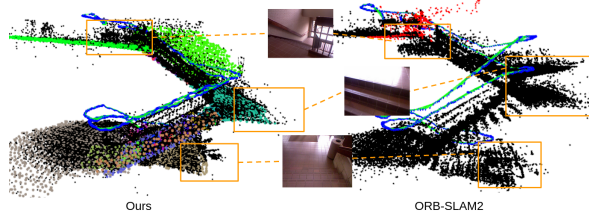


Fig. 5. Qualitative results of sparse reconstruction and trajectory between the proposed method and ORB-SLAM2 in the TAMU dataset.

that in our experiments, the initial search radius for selecting neighbors for triangulation is set to $5m$ and the multiplier is set as 5 to modify the final search radius to adapt to different point densities on the plane regions.

VI. EXPERIMENTS

We evaluate the proposed SLAM system on two well known public datasets, the ICL-NUIM [34] and TUM RGB-D [35] benchmarks, comparing its performance with other state-of-the-art methods such as ORB-SLAM2 [7], PS-SLAM [6] that are feature-based methods, but removed the global bundle adjustment modules in the following experiments. Methods based on the MW assumption such as LPVO [11] and L-SLAM [12]. DVO-SLAM [18] is a direct method and InfiniTAM [32] uses a GPU for real-time tracking and mapping based on RGB and depth images. Additionally, we provide the reconstruction accuracy of our

time	Feat. extr.	Rotat.	Transla	Refinement	Total
Median	19.9	2.1	4.8	13.0	42.5
Mean	20.5	3.0	5.4	13.1	43.7
Std.	3.6	0.4	2.8	4.8	9.4

TABLE II

MEASURED TRACKING TIMES (MS) ON THE TUM RGB-D SEQUENCES

Sequence	RGB-D	ElasticFu	InfiniTAM	SPFu	Ours
kt0	4.4	0.7	1.3	0.7	0.4
kt1	3.2	0.7	1.1	0.9	0.6
kt2	3.1	0.8	0.1	1.1	0.6
kt3	16.7	2.8	2.8	1.0	0.8

TABLE III

RMSE RECONSTRUCTION ERROR (CM) ON THE ICL-NUIM DATASET IN CENTIMETERS.

reconstructed model on the ICL-NUIM dataset and compare it with other popular methods for dense reconstruction. Lastly, to demonstrate that our system is robust over time, we also test on a sequence from the TAMU [3] dataset containing long sequences covering a large indoor area. All experiments are carried out with an Intel Core i7-8700 CPU (with @3.20GHz) and without any use of GPU. The ICL-NUIM dataset [34] provides synthetic scenes for two indoor environments, one living room and one office room scenario. These scenes contain large areas of low textured surfaces such as walls, ceilings, floors, etc. There are four sequences for each scene. We evaluate our method on all sequences.

A. ICL-NUIM RGB-D Dataset

Table I shows that our method obtains the best performance on three out of the eight sequences, based on the translation RMSE (ATE). InfiniTAM also performs well on lr-kt1, lr-kt2 and of-kt3 sequences, but the method also loses tracking in other sequences. As the dataset contains large structured areas, the Manhattan-based methods LPVO and L-SLAM are able to get a good estimate of the orientation and provide good results throughout. However, they usually

Sequences	Ours	ORB-SLAM2	length
Corridor-A	1.62	3.13	88
Stair-A	0.94	1.44	66
Entry-Hall	1.33	2.22	80

TABLE IV

COMPARISON OF THE ACCUMULATED DRIFT (M) IN DIFFERENT LARGE SCALE SEQUENCES.

need two planes, or alternatively, one plane and a vanishing direction to be visible at all times to estimate a good Manhattan frame. As shown in Figure 6, there are several

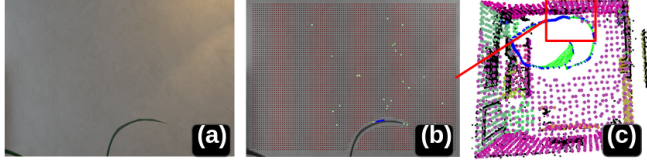


Fig. 6. Results in lr-k3. (a) input image; (b) point-line features and the segmented plane; (c) reconstructed 3D map and trajectory.

challenging scenes in lr-kt3, where only a white wall and two leaves from a plant are captured when the camera is close to the wall. In this situation, OPVO and L-SLAM are unable to yield a good performance. When a bad initial pose is obtained in our system due to the scene not being a rigid MW, the refinement step based on point-line-plane features allows us to recover the pose nevertheless, while L-SLAM ignores optimizing rotation in the LKF module. Moreover, while DVO, being a dense method, may struggle because of the large areas of walls, floor etc. not containing enough gradient for the photometric error, ORB-SLAM2 and PS-SLAM perform well, as both environments contain sufficient ORB features extracted from furniture, objects etc. As our method takes advantage of all geometric elements, it is able to perform robustly in most sequences. In addition, Figure 4 shows the relative pose error for ORB-SLAM, PS-SLAM and our method. Our method obtains notably better results than the other two in relative translation and rotation. Especially the rotation error is much lower for our method, due to the use of the decoupled MW rotation estimation.

B. TUM RGB-D Dataset

The TUM RGB-D benchmark [35] is one of the most popular datasets for RGB-D SLAM systems, which provides indoor sequences under different texture and structure conditions. This allows us to separately test sequences which have structure, texture or both. In order to evaluate our method in challenging environments, we select four structured image sequences, the first three with low texture and the last one with a large scale environment. As all sequences listed in Table I have structure, but the large-cabinet sequence is not a rigid Manhattan scenario. Manhattan-based methods are able to provide good pose estimates on snot-far sequence, but the results degenerate in large-cabinet and cabinet sequences. The first two sequences include the same environment consisting of multiple non-textured planes. Here ORB-SLAM2

is not able to find enough point correspondences along the sequence and loses tracking. Our method, which additionally uses lines and planes for translation estimation, achieves better results. As shown in Figure 4, cabinet and large-cabinet are challenging sequences because of several low-texture frames. Our method's tracking strategy limits the relative rotation error to under 2 degrees, which is better than ORB-SLAM2 and PS-SLAM. The statistics of the time spent for each operation are shown in Table II, where we use different CPU threads to deal with points, lines and planes in the feature extraction and refinement modules.

C. Large scale sequence

The TAMU dataset [36] provides large-scale indoor sequences (constant lighting). While it does not provide ground-truth camera poses, the start and end point are the same, which can be used to evaluate the overall drift by computing the final position errors. As shown in Figure 5, the trajectory in the sequence Stair-C is a loop between two floors, where the improvement of our method over the whole trajectory length is 34.7% in drift compared to ORB-SLAM2. Similar situations can also be found in Corridor-A and Entry-Hall. More qualitative results are provided in the supplementary material.

D. Reconstruction Accuracy

We reconstruct models from ICL-NUIM and compare the results with state-of-the-art mapping methods, as shown in Table III. The accuracy of the reconstruction results is defined as the mean difference between the predicted model and the ground-truth model [34]. We compare the proposed mapping module against RGB-D SLAM [37], ElasticFusion [9], InfiniTAM [38], and SuperpixelFusion [10].

The SuperpixelFusion method is constrained by using ORB-SLAM for pose estimation, whereas our method also works well in low-textured environments. InfiniTAM obtains the best results in kt2, but shows worse performance on the kt0 and kt3 sequences, potential due to the large low-textured regions. ElasticFusion shows a similar behavior. Our method reconstructs more accurate maps than the others, but InfiniTAM and ElasticFusion provide more complete models than our map since we ignore small objects even though features based on points, lines and planes cover most of the pixels. Remarkably, all fusion methods, except for SuperpixelFusion and ours, rely on GPU based acceleration.

VII. CONCLUSIONS

We have proposed a RGB-D SLAM system based on points, lines and planes. Using the MW assumption for rotation estimation, and point, line and plane features for translation estimation, we achieve state-of-the-art performance. Also, a novel instance-wise meshing approach can reconstruct planar regions in the environment efficiently. The resulting dense map allows for interactions with the environment in robotic and AR/VR applications. In the future we would like to extend the planar reconstruction with a meshing of the non-planar parts in the environment to allow the complete reconstruction of more complex scenes.

REFERENCES

- [1] Y. Li, N. Brasch, Y. Wang, N. Navab, and F. Tombari, "Structure-slam: Low-drift Monocular SLAM in Indoor Environments," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2020.
- [2] R. Gomezjeda, F. Moreno, D. Scaramuzza, and J. G. Jimenez, "PL-SLAM: a Stereo SLAM System Through the Combination of Points and Line Segments," *IEEE Trans. Robot.*, 2019.
- [3] Y. Lu and D. Song, "Robust RGB-D Odometry Using Point and Line Features," in *Proc. IEEE Int. Conf. Comput. Vision*, 2015.
- [4] A. Pumarola, A. Vakhitov, A. Agudo, A. Sanfeliu, and F. Morenonoguer, "PL-SLAM: Real-time Monocular Visual SLAM with Points and Lines," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2017.
- [5] P. F. Proença and Y. Gao, "Probabilistic RGB-D odometry based on points, lines and planes under depth uncertainty," *Robotics and Autonomous Syst.*, vol. 104, pp. 25–39, 2018.
- [6] X. Zhang, W. Wang, X. Qi, Z. Liao, and R. Wei, "Point-plane SLAM using supposed planes for indoor environments," *Sensors*, vol. 19, no. 17, p. 3795, 2019.
- [7] R. Murartal and J. D. Tardos, "ORB-SLAM2: An Open-Source SLAM System for Monocular, Stereo, and RGB-D Cameras," *IEEE Trans. Robot.*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [8] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, "KinectFusion: Real-time Dense Surface Mapping and Tracking," in *IEEE Int. Symposium on Mixed and Augmented reality*, 2011.
- [9] T. Whelan, S. Leutenegger, R. Salas-Moreno, B. Glocker, and A. Davison, "ElasticFusion: dense SLAM without a pose graph," in *Robotics: Science and Syst.*, 2015.
- [10] K. Wang, F. Gao, and S. Shen, "Real-time Scalable Dense Surfel Mapping," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2019.
- [11] P. Kim, B. Coltin, and H. J. Kim, "Low-drift Visual Odometry in Structured Environments by Decoupling Rotational and Translational Motion," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2018.
- [12] P. Kim, B. Coltin, and H. Jin Kim, "Linear RGB-D SLAM for Planar Environments," in *Proc. Springer Eur. Conf. Comput. Vision*, 2018.
- [13] C. Raposo, M. Lourenço, M. Antunes, and J. P. Barreto, "Plane-based odometry using an RGB-D camera," in *Proc. Brit. Mach. Vision Conf.*, 2013.
- [14] Y. Taguchi, Y.-D. Jian, S. Ramalingam, and C. Feng, "Point-plane SLAM for hand-held 3D sensors," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2013.
- [15] R. F. Salas-Moreno, B. Glocker, P. H. Kelly, and A. J. Davison, "Dense Planar SLAM," in *IEEE Int. Symposium on Mixed and Augmented reality*, 2014.
- [16] P.-H. Le and J. Košečka, "Dense piecewise planar RGB-D SLAM for indoor environments," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2017.
- [17] F. Nardi, B. Della Corte, and G. Grisetti, "Unified representation and registration of heterogeneous sets of geometric primitives," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 625–632, 2019.
- [18] C. Kerl, J. Sturm, and D. Cremers, "Dense visual SLAM for RGB-D cameras," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2013.
- [19] L. Ma, C. Kerl, J. Stückler, and D. Cremers, "CPA-SLAM: Consistent Plane-model Alignment for Direct RGB-D SLAM," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2016.
- [20] T. Schops, T. Sattler, and M. Pollefeys, "Bad SLAM: bundle adjusted direct RGB-D SLAM," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2019.
- [21] T. Schöps, T. Sattler, and M. Pollefeys, "Surfel-meshing: online surfel-based mesh reconstruction," *IEEE Trans. on Pattern Analysis and Machine Intell.*, 2019.
- [22] S. Schreiberhuber, J. Prankl, T. Patten, and M. Vincze, "ScalableFusion: High-resolution Mesh-based Real-time 3D Reconstruction," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2019.
- [23] J. Straub, N. Bhandari, J. J. Leonard, and J. W. Fisher, "Real-time Manhattan World Rotation Estimation in 3D," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2015.
- [24] Y. Zhou, L. Kneip, C. Rodriguez, and H. Li, "Divide and Conquer: Efficient Density-based Tracking of 3D Sensors in Manhattan Worlds," in *Proc. Asian Conf. Comput. Vision*, 2016.
- [25] K. Joo, T.-H. Oh, J. Kim, and I. So Kweon, "Globally Optimal Manhattan Frame Estimation in Real-time," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2016.
- [26] P. Kim, B. Coltin, and H. J. Kim, "Visual Odometry with Drift-free Rotation Estimation Using Indoor Scene Regularities," in *Proc. Brit. Mach. Vision Conf.*, 2017.
- [27] E. Rublee, V. Rabaud, K. Konolige, and G. R. Bradski, "ORB: An efficient alternative to SIFT or SURF," in *Proc. IEEE Int. Conf. Comput. Vision*, 2011.
- [28] R. G. Von Gioi, J. Jakubowicz, J.-M. Morel, and G. Randall, "LSD: A Fast Line Segment Detector with a False Detection Control," *IEEE Trans. on Pattern Analysis and Machine Intell.*, vol. 32, no. 4, pp. 722–732, 2010.
- [29] L. Zhang and R. Koch, "An Efficient and Robust Line Segment Matching Approach Based on LBD Descriptor and Pairwise Geometric Consistency," *Elsevier Journal of Visual Communication and Image Representation*, vol. 24, no. 7, pp. 794–805, 2013.
- [30] A. J. Trevor, S. Gedikli, R. B. Rusu, and H. I. Christensen, "Efficient Organized Point Cloud Segmentation with Connected Components," *Semantic Perception Mapping and Exploration*, 2013.
- [31] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "ORB-SLAM: a Versatile and Accurate Monocular SLAM System," *IEEE Trans. Robot.*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [32] V. A. Prisacariu, O. Köhler, S. Golodetz, M. Sapienza, T. Cavallari, P. H. Torr, and D. W. Murray, "Infinitam v3: A Framework for Large-scale 3D Reconstruction with Loop Closure," *arXiv preprint arXiv:1708.00783*, 2017.
- [33] Z. C. Marton, R. B. Rusu, and M. Beetz, "On Fast Surface Reconstruction Methods for Large and Noisy Datasets," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2009.
- [34] A. Handa, T. Whelan, J. McDonald, and A. Davison, "A Benchmark for RGB-D Visual Odometry, 3D Reconstruction and SLAM," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2014.
- [35] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A Benchmark for the Evaluation of RGB-D SLAM Systems," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2012.
- [36] Y. Lu and D. Song, "Robustness to Lighting Variations: An RGB-D Indoor Visual Odometry Using Line Segments," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2015.
- [37] F. Endres, J. Hess, N. Engelhard, J. Sturm, D. Cremers, and W. Burgard, "An Evaluation of the RGB-D SLAM System," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2012.
- [38] O. Köhler, V. A. Prisacariu, and D. W. Murray, "Real-time Large-scale Dense 3D Reconstruction with Loop Closure," in *Proc. Springer Eur. Conf. Comput. Vision*, 2016.