

Visual SLAM with Graph-Cut Optimized Multi-Plane Reconstruction

Fangwen Shu *

Yaxu Xie

Jason Rambach

Alain Pagani †

Didier Stricker

DFKI - German Research Center for Artificial Intelligence

ABSTRACT

This paper presents a semantic planar SLAM system that improves pose estimation and mapping using cues from an instance planar segmentation network. While the mainstream approaches are using RGB-D sensors, employing a monocular camera with such a system still faces challenges such as robust data association and precise geometric model fitting. In the majority of existing work, geometric model estimation problems such as homography estimation and piece-wise planar reconstruction (PPR) are usually solved by standard (greedy) RANSAC separately and sequentially. However, setting the inlier-outlier threshold is difficult in absence of information about the scene (i.e. the scale). In this work, we revisit these problems and argue that two mentioned geometric models (homographies/3D planes) can be solved by minimizing an energy function that exploits the spatial coherence, i.e. with graph-cut optimization, which also tackles the practical issue when the output of a trained CNN is inaccurate. Moreover, we propose an adaptive parameter setting strategy based on our experiments, and report a comprehensive evaluation on various open-source datasets.

Index Terms: Computing methodologies—Artificial intelligence—Computer vision—Computer vision tasks—Vision for robotics

1 INTRODUCTION

Semantic planar SLAM has gained much attention in the last decade, especially for virtual reality (VR) systems and augmented reality (AR) applications. Although there has been intensive research on this topic, most of the current methods still focus on RGB-D sensors [9, 11, 21, 24] with plane primitives extracted from depth images. While monocular methods [19, 25, 27] still face several challenges and difficulties, such as low-texture scenes, dynamic foregrounds, pure rotation of the camera, various baseline between frames, and scale drift, where plane primitives can only be extracted from limited 3D information obtained. Existing methods either build upon indirect SLAM [18] or direct SLAM [6], and both are subjected to the challenges mentioned before.

In this work, we argue that data association and geometric model fitting problems are usually not tackled efficiently in monocular SLAM systems, i.e. establishing feature matches of multi-plane between frames taken from different viewpoints (under small or large baseline), or from the same viewpoint (under pure rotation), with homographies estimated and decomposed. Thereafter, in order to localize the camera relatively, a plausible homography matrix is usually verified by triangulation (via positive depth validation) and minimizing symmetric transfer error (STE) between image pairs. However, the map scale is not observable purely from relative pose estimation. At the same time, 3D planes can only be fitted from sets of noisy and sparse point clouds triangulated by the monocular SLAM. Thus, to tackle the problems especially for monocular systems, we first integrate a real-time instance planar segmentation network into a feature-based SLAM system. Then we propose to solve the multi-model fitting problem in a sequential RANSAC fashion but with a fast graph-cut optimized proposal engine [2]. Finally,

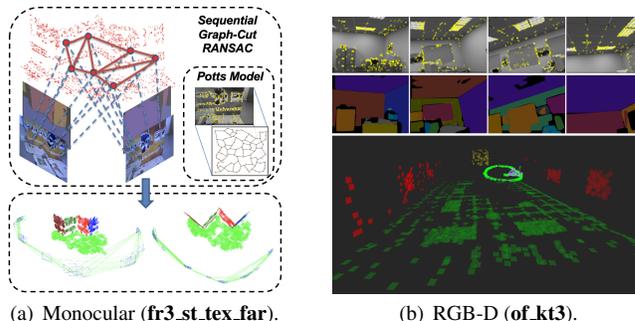


Figure 1: **We propose using Sequential Graph-Cut RANSAC (Algo. 1) with feature-based SLAM for robust piece-wise planar reconstruction (PPR).** Here, the illustrated light-weight, planar, patch-surfels semantic map is reconstructed from sparse and noisy point clouds. Every distinguished color indicates a different plane.

we present quantitative and qualitative results from our semantic SLAM system, which shows the proposed method can be applied on any feature-based monocular or RGB-D SLAM system without a significant algorithm adaptation. To summarize, we propose the following contributions:

- We introduce an energy-based geometric model fitting method, i.e. sequential RANSAC with graph-cut optimization, into a feature-based planar SLAM system, which implicitly considers SLAM as optimizing geometric multi-model estimation of different types.
- We propose a SLAM building block that tightly integrates the energy-based method mentioned above and a state-of-the-art convolutional neural network (CNN) of instance planar segmentation. Thus, we do not take any output from CNN as a noise-free "sensor" measurement, but further optimize it within the SLAM workflow, which boosts the performance of tracking and optimization.
- We conduct exhaustive experiments and report a comprehensive evaluation on various indoor datasets.

2 RELATED WORK

RGB-D Planar SLAM. There is a large body of recent and ongoing research on planar SLAM using 3D sensors. A common scenario of using semantic planar cues in SLAM is adding geometric regularization regarding different landmarks and optimize the geometric structure jointly, such as with Manhattan World (MW) assumption [28]. An early work [24] presented a SLAM system for a hand-held 3D sensor, where the authors argued that it is possible to register 3D data in two different coordinates using any combination of point and plane primitives. This was followed by the works [9, 11, 21, 29] which tackle the problem similarly by extracting plane primitives from depth image, registering planes across different views, and optimize the poses of keyframes and landmarks (both point and plane) in Bundle Adjustment (BA).

Monocular Planar SLAM. Even if the MW assumption is a good constraint for indoor SLAM, it is difficult to enforce it in monocular methods because only limited 3D information can be obtained. Therefore, this strong assumption is generally not used in the case

*e-mail: {first_name}.{last_name}@dfki.de

†**Acknowledgment:** This research is partially funded by the German BMBF project MOVEON (01IS20077) and SocialWear (01IW20002).

of a monocular sensor. In this work, we focus more on developing a robust multi-plane estimator and refinement strategy. For example, π Match [20] employs PEARL (an energy-based geometric multi-model fitting method [10]) for piece-wise planar reconstruction in a novel two-view Structure-from-Motion (SfM) workflow. A recent monocular SLAM framework [27] uses the high-level object and plane landmarks, leading to a dense, compact, and semantically meaningful map compared to the classic feature-based SLAM. To track in a low-texture environment, Structure-SLAM [14] employs not only planes but also predicted normals and lines to calculate drift-free rotation. More recently, [25] proposed a Winner-Takes-All (WTA) RANSAC-based relative camera pose estimation under multiple planar structures with the help of superpixels segmentation.

Exploiting the Spatial Coherence in RANSAC. This kind of method is usually formalized as a binary labeling problem for geometric multi-model fitting, which we will discuss in detail in Sec. 3. First with PEARL (Propose Expand and Re-estimate Labels) [10], then with more advanced methods such as Graph-Cut RANSAC [2] and Progressive-X [3], the spatial coherence is exploited in the local optimization step to find local structures accurately, which consider that geometric data often form spatially coherent structures described by, e.g., the Potts model [4]. Those methods are suitable for our problems in both 2D and 3D space, where homographies and piece-wise planes should not have spatial overlapping.

Plane Detection and Reconstruction via CNN. Initially PlaneNet [16] presented an end-to-end neural architecture for piece-wise planar reconstruction from a single RGB image. However, it suffers from limitations such as missing small surfaces and requiring the maximum number of planes in a single image as prior. PlaneRCNN [15] addressed these issues and proposed the first detection-based neural network for piece-wise planar reconstruction, which jointly refines all the segmentation masks with warping loss function to enforce the consistency with a nearby view during training. It is able to detect small planar surfaces, but fails to reach a real-time frame rate. More recently, PlaneSegNet [26] was proposed as a real-time one-stage instance plane segmentation network that achieves significantly higher frame rates and comparable segmentation accuracy against the two-stage methods mentioned before.

3 METHOD

We first introduce the standard sequential RANSAC pipeline for geometric model fitting (homography or plane structure in this work), with semantic cues as input with image sequences. However, we would like to cope with possible misclassification from the instance segmentation network, and therefore we do not simply use a standard RANSAC-like plane fitting algorithm for each detected planar segment. Instead, we propose a more robust sequential pipeline using a locally optimized RANSAC alternating graph-cut and model re-fitting in the inner local optimization step (Algo. 1) to adapt automatically to inaccurate instance segmentation and noise. Finally, we discuss how it is integrated into a feature-based SLAM framework (as shown in Fig. 2) as a robust geometric multi-model fitting strategy in this work. The mathematical notation used in this section is adopted from [2, 10].

3.1 Geometric Model Fitting via RANSAC

Standard RANSAC [7] is a well-known method for dealing with a large number of outliers when data supports only one model (e.g. 2D line fitting). The implied unary energy: $E_{\{0,1\}}(L) = \sum_p \|L_p\|_{\{0,1\}}$ counts inliers for the target model using 0-1 measure, thus can be reformulated as binary labeling problem [10], where parameter vector θ of the model with the largest number of inliers

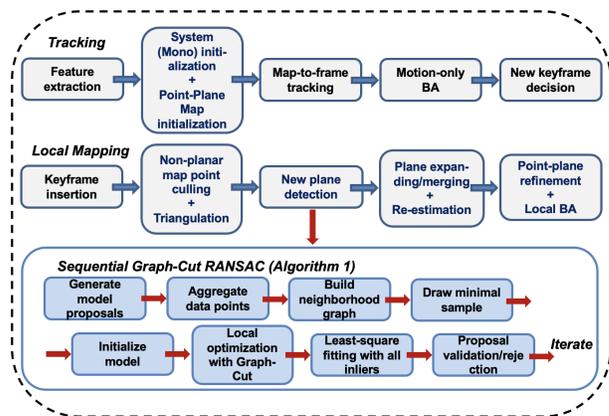


Figure 2: The overall structure/workflow of the monocular planar SLAM system used in this work.

within some threshold ϵ is estimated:

$$\|L_p\|_{\{0,1\}} = \begin{cases} 0 & \text{if } (L_p = 1 \wedge \text{dist}(p, \theta) < \epsilon) \vee \\ & (L_p = 0 \wedge \text{dist}(p, \theta) \geq \epsilon) \\ 1 & \text{otherwise.} \end{cases}$$

where $\|L_p\|_{\{0,1\}}$ is the geometric error measure, with $\text{dist}(p, \theta)$ the distance function, i.e. the Euclidean distance between data point $p \in P$ and the estimated model. Parameter $L \in \{0, 1\}^{|P|}$ is the labeling, $|P|$ is the number of data points. $L_p \in L$ is the label of the point p . Here, the unary energy penalize nothing when p is labeled as inlier (close to the model) or p is labeled as outlier (far from the model).

Sequential RANSAC. As a variant [30] of RANSAC, sequential RANSAC detects model instances one after another, with the inliers of the detected instance removed from data point set P . The drawback of this approach is that the inliers are usually assigned to the model instance which has the most support instead of the actual best instance. When we apply it with semantic priors, the proposed model's inliers are already clustered from the data point set, thus it is followed by a simple validation of the model. The result is, obviously, dependent on the generalizability of a trained CNN and the user-defined thresholds that are difficult to fine-tune.

3.2 Sequential Model Fitting with Spatial Coherence

The geometric multi-model fitting problem is usually formulated as an optimal labeling problem, where the binary energy $E(L)$ can be extended with an additional term indicating the label count penalty (label smoothness) [13] and a term indicating the spatial regularity [10]. The existing methods usually consider the general case when the data supports some unknown number of models, and solving $E(L)$ over labeling $L = \{L_p | p \in P\}$ which describes the overall quality of the solution. However, in our case, the number of models is known from CNN, thus making the problem slightly different.

Energy Formulation. In this work, we assume the number of models is known, so there is no label count penalty added. We, therefore, formulate geometric fitting problems simply by optimizing the energy $E(L)$ over K different models (L_1, L_2, \dots, L_K) , in a sequential way. The implicit assignment of inliers to models becomes trivial in our case because of an available segmentation prior, which means generating a large number of proposed labels (models) done in PEARL [10] is not needed anymore, making the algorithm very fast. In this case, for each model (labeling) proposal L , the assigned data point p will be labeled as inlier or outlier after a so-far-the-best-model is found. As both inliers and outliers of the model should be

Algorithm 1: Sequential Graph-Cut RANSAC.

Input : M_{seg} – segmentation mask;
 P – extracted feature points or matched correspondences;
 N – max. iteration number of outer RANSAC loop;
 N_{GC} – max. iteration number of inner GC-RANSAC loop;
 ϵ_d – point-to-model distance threshold;
 ϵ_L – residual threshold of a model;

Output : θ – parameters of model instances;

```

1 for  $p \in P$  do
2    $L_p \leftarrow$  Generate model (labeling) proposals using  $M_{seg}$ ;
3    $V_p \leftarrow$  Aggregate the data point;
4 end
5 for  $L_p \in L$  do
6    $e_L^* \leftarrow$  Initialize residual with the max. numeric value;
7    $\mathcal{N} \leftarrow$  Build neighborhood-graph from  $V_p$ ;
8   for  $i = 1 \rightarrow N$  do
9     for  $k = 1 \rightarrow N_{GC}$  do
10       $S_k \leftarrow$  Draw a minimal sample from  $V_p$ ;
11       $\theta_k \leftarrow$  Estimate model parameters using  $S_k$ ;
12       $w_k \leftarrow$  Find the inliers of  $\theta_k$  using  $\epsilon_d$ ;
13       $e_{L_k} \leftarrow$  Calculate the residuals;
14      if  $e_L^* > e_{L_k}$  then
15         $e_L^*, \theta_k^*, w_k^* \leftarrow e_{L_k}, \theta_k, w_k$ ;
16        if Local opt. (refer to Algo. 2 [2]) then
17           $w_{LO}^*, changed \leftarrow w_k, 1$ ;
18          while changed do
19             $G \leftarrow$  Build the problem graph (refer to
                Algo. 3 [2]);
20             $L_{LO} \leftarrow$  Apply graph-cut to  $G$ ;
21             $I_{7m} \leftarrow$  Select a  $7m$ -sized random inlier
                set;
22             $\theta_{LO} \leftarrow$  Fit a model using  $I_{7m}$ ;
23             $w_{LO} \leftarrow$  Compute the support of  $\theta_{LO}$ ;
24             $changed \leftarrow 0$ ;
25            if  $w_{LO} > w_{LO}^*$  then
26               $\theta_{LO}^*, L_{LO}^*, w_{LO}^*, changed \leftarrow$ 
                 $\theta_{LO}, L_{LO}, w_{LO}, 1$ ;
27            end
28          end
29           $L_p^*, \theta_k^*, w_k^* \leftarrow L_{LO}^*, \theta_{LO}^*, w_{LO}^*$ ;
30        end
31      end
32       $e_L^*, \theta^* \leftarrow$  Least squares fitting (SVD) using  $w_k^*$ ;
33    end
34    if  $e_L^* < \epsilon_L$  then
35      Early break;
36    end
37  end
38 end

```

spatially coherent, which means a point near an outlier (resp. inlier) is more likely to be an outlier (resp. inlier). In order to take this into account, we propose to use a graph-cut algorithm to further optimize the inlier-outlier assignment. The proposed energy:

$$E(L) = \sum_p \|L_p\| + \lambda \cdot \sum_{(p,q) \in \mathcal{N}} w_{pq} \cdot \delta(L_p \neq L_q) \quad (1)$$

where the first term indicates the geometric error measure between data point and corresponding model, and the second term indicates the spatial regularization which penalizes neighbors with different labels in the graph. \mathcal{N} indicates edges in the near-neighbor graph constructed from data point set (e.g. the Potts model in Fig. 1). $\delta(\cdot)$ is 1 if the specified condition inside parenthesis holds, and 0 otherwise. Weights w_{pq} set discontinuity penalties for each pair of neighboring data points. λ is a parameter balancing the two terms.

The binary labeling energy minimization with additional spatial

regularity terms (Eq. 1) can be solved efficiently and globally via the graph-cut algorithm. We adopted the idea of conducting the graph-cut algorithm in the local optimization (LO) step which is applied when a so-far-the-best model is found [2]. As the local optimization [5] assumes not all all-inlier samples are "good", it is perfectly aligned with our problem in this work, where the instance segmentation clusters the data samples in the first stage is prone to be inaccurate. More importantly, it is real-time feasible to apply the graph-cut in the LO step within just a few iterations, as the local optimization step converges very fast when it takes spatial proximity into account. The proposed *Sequential Graph-Cut RANSAC* is presented in Algo. 1. The whole pipeline can be considered as several steps: (1) generate model proposals (labeling L_1, L_2, \dots, L_k) based on segmentation mask and assign the data points to the model; (2) estimate each model proposal sequentially, where the construction of problem graph G (line 19) within the local optimization is used to build energy minimization of Eq. 1; (3) thereafter the graph-cut is applied to G determining the optimal labeling L ; (4) model parameters are updated according to, not only the number of the support inliers w as done in [2], but also a residual threshold ϵ_L defined within the SLAM system used in this work. Please note that the used parameters and thresholds will be further explained and detailed in the experiments (Sec. 4.2).

3.3 Visual SLAM Framework

System Initialization and Map Initialization. For monocular SLAM, we establish the proposed Algo. 1 within the initialization step, where the Homography matrix and Fundamental matrix are calculated in parallel as done in [18]. We use symmetric transfer error (STE) for the geometric error measure $\|H_p\|$ between matched feature points $p = (p_{ref}, p_{cur})$. The initial solution for the non-linear minimization is found by using the Normalized Direct Linear Transform (NDLT) with the minimal 4 correspondences. Then we apply the energy minimization (Eq. 1) for homographies:

$$E(\mathbf{H}) = \sum_p \|H_p\| + \lambda \cdot \sum_{(p,q) \in \mathcal{N}} \delta(H_p \neq H_q) \quad (2)$$

where $\mathbf{H} = \{H|p \in P\}$ is the assignment of models to feature points p in the reference frame, the neighborhood system \mathcal{N} is based on a grid-neighborhood construction on image space and the minimum samples (4 correspondences) are sampled by progressive-NAPSAC sampler [1] within that image grid.

The homography with the most inliers is used to calculate the score S_H [18] and initialize the map. The fundamental matrix is calculated using the default implementation. Thereafter the initial map is scaled by setting the median of the inverse of the depth as 1 before tracking the next frame. After re-scaling the map, several planes can be fitted from the 3D point cloud. To this aim, we apply the energy minimization (Eq. 1) again within Algo. 1 for piece-wise planar reconstruction:

$$E(\mathbf{\Pi}) = \sum_v \|\Pi_v\| + \lambda' \cdot \sum_{(u,v) \in \mathcal{N}'} \delta(\Pi_u \neq \Pi_v) \quad (3)$$

where $\mathbf{\Pi} = \{\Pi|v \in V\}$ is the assignment of plane models to map point, and V indicates the set of 3D vertices. Here, we use the distance between a 3D point and a plane: $dist(v, \Pi) = \frac{|\mathbf{n}^T \cdot \mathbf{v} + d|}{\|\mathbf{n}\|}$ as geometric error measure in the first term $\|\Pi_v\|$, with the plane represented as $(\mathbf{n}^T, d)^T$, where \mathbf{n} is the plane normal and d is the distance to the world origin. The neighborhood system \mathcal{N}' is constructed using Fast Approximate Nearest Neighbors algorithm [17] according to a predefined sphere radius r as the 3D grid is unknown, and the minimum samples (3 points) are sampled uniformly.

Local Mapping with Plane Expanded and Re-estimated. The existing feature-based method (i.e. ORB-SLAM2 [18]) focuses on

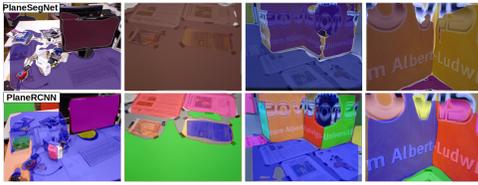


Figure 3: **Comparison of the segmentation results** of PlaneSegNet [26] and PlaneRCNN [15] on dataset TUM RGB-D [22].

utilizing as much point data as possible. When a new keyframe is inserted, however, we remove (local) map point which is assumed being associated with a plane but its distance $dist(v, \Pi)$ is bigger than a threshold ϵ_d . This step is conducted before triangulation and local BA, thus will not influence the stability of the system. Right after the new point landmark is triangulated (still before local BA), we conduct Algo. 1 for piece-wise planar reconstruction. However, detecting outliers of the geometric model is somewhat heuristic in this work. To avoid the presence of planes with weak support in terms of number of points, we consider planes with a low number of 3D points as "weak" planes and we only keep high-quality planes in the map. The weak planes can be later merged into other plane instances or removed if they cannot be expanded.

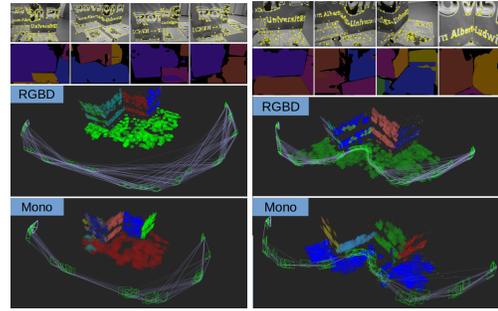
A map refiner is running in a loop within the local mapping thread, checking all the 3D plane instances and trying to merge the closed planes. Two planes are merged if the following two conditions are met, first they should have nearly parallel normals: $|\cos(\theta)| = \frac{|n_i n_j|}{\|n_i\| \|n_j\|} > T_\theta$ ($0 < T_\theta < 1$) and second, they should be geometrically close to each other: $|\frac{d_i}{\|d_i\|} - \frac{d_j}{\|d_j\|}| < T_d$. The new plane equation is updated in a RANSAC loop with 60% randomly sampled associated point landmark. After that, all associated point landmarks will be projected on the plane by minimizing the point-plane distance via: $\hat{v} = v - dist(v, \Pi_v) \cdot \frac{n}{\|n\|}$.

Optimization. Different Bundle Adjustment (BA) are considered: (1) Motion-only BA for map-to-frame tracking. (2) Local BA on keyframe will be conducted after new plane detected, existing plane merged/expanded and re-estimated. Notice for the local BA we treat the optimization of the structure and the motion in a de-coupled way. It means the reconstructed plane is used to structure the map first, thereafter the map is used in the local BA for joint optimization. (3) Global BA only happens after loop closure. While a final refinement could be conducted with cost function [12]:

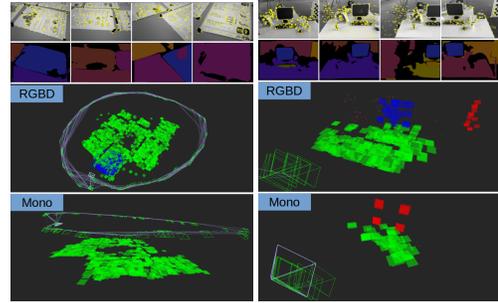
$$\operatorname{argmin}_{C_i, v_j, \Pi_k} \sum_j \left\{ \sum_i \|p_{ij} - Q(C_i, v_j)\| + \sum_k dist(v_j, \Pi_k) \right\} \quad (4)$$

where the first term indicates the standard reprojection error and $Q(\cdot)$ is the camera projection function. The second term indicates the point-to-plane distance. The 6-DOF camera pose is represented as Lie Algebra $C \in \mathfrak{se}(3)$. i, j, k are the number of camera views, 3D points and planes, respectively. Notice that the plane needs to be parameterized as minimal representation such as spherical coordinates: $\Pi = (\phi = \arctan(\frac{n_y}{n_x}), \psi = \arcsin(n_z), d)$, where ϕ and ψ are the azimuth and elevation angles of the plane normal.

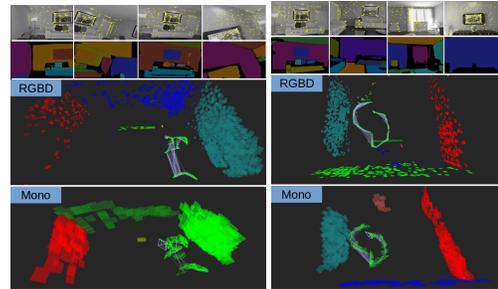
However, under monocular setting, this cost function is subject to outlier planes (e.g. data sequences such as fr1_desk/fr2_desk, explained in Sec. 4.2) and can become a large scale non-linear optimization (e.g. any corridor scenario or large environment). In this work, for the fair comparison with other semantic SLAM methods, we report quantitative results (in Table 1) without the final global optimization, and focus on the local mapping of PPR which shows the direct impact of the PPR on the accuracy of estimated trajectory.



(a) fr3_st_tex_far. (b) fr3_st_tex_near.



(c) fr3_nst_tex_near. (d) fr2_xyz.



(e) lr_kt0. (f) lr_kt2.

Figure 4: **The light-weight semantic map (points and planar patch-surfels, best view zoom-in)** constructed on selected sequences of dataset TUM RGB-D [22] and ICL-NUIM [8].

4 EXPERIMENTS AND RESULTS

We base our experiments on datasets TUM RGB-D [22] and ICL-NUIM [8]. First, we discuss the chosen instance plane segmentation CNN, especially highlighting the problematic failure cases and explain how we treat it as noisy "sensor" measurement. Then we compare the performance of our SLAM framework with the proposed multi-plane estimator to some recent vSLAM methods.

4.1 Plane Detection via CNN

We employ PlaneSegNet [26] as our plane detector, because, as a global-mask-based instance segmentation method, it provides segmentation masks with higher resolution and better completeness than local-mask based method like PlaneRCNN [15]. More important, PlaneSegNet is capable to run at a real-time frame rate (over 30Hz), while PlaneRCNN is only able to run at less than 5Hz with the same hardware (NVIDIA GTX 1080 Ti).

The instance segmentation is considered as a prior information for piece-wise planar reconstruction. However, as shown in Fig. 3, global-mask based instance segmentation method (i.e. PlaneSegNet) suffers from feature leakage, and sometimes cannot distinguish dif-

| Dataset & Sequences | Monocular | | | | RGB-D | | | | |
|-----------------------|-------------|----------------|-----------------|---------------------|-------------|----------------|-----------------|---------------------|--------------|
| | Ours | ORB-SLAM2 [18] | Open-VSLAM [23] | Structure-SLAM [14] | Ours | ORB-SLAM2 [18] | Open-VSLAM [23] | Manhattan SLAM [28] | SP-SLAM [29] |
| TUM RGB-D [22] | | | | | | | | | |
| fr1_xyz | 0.93 | 0.99 | 1.12 | - | 1.02 | 1.20 | 1.25 | 1.00 | 0.93 |
| fr1_floor | 1.72 | 2.99 | 1.91 | - | 1.64 | 2.50 | 1.79 | 2.50 | - |
| fr1_desk | 1.91 | 1.41 | 1.73 | - | 1.81 | 1.73 | 1.86 | 2.70 | 1.43 |
| fr2_xyz | 0.24 | 0.27 | 0.26 | - | 1.73 | 0.37 | 1.71 | 0.80 | - |
| fr2_desk | 1.28 | 1.14 | 0.90 | - | 7.66 | 0.85 | 7.85 | 3.70 | - |
| fr3_st.tex_far | 0.84 | 0.95 | 1.09 | 1.40 | 1.08 | 1.19 | 1.10 | 2.20 | 0.97 |
| fr3_st.tex_near | 1.24 | 1.29 | 1.28 | 1.40 | 0.83 | 1.21 | 0.91 | 1.20 | 0.84 |
| fr3_nst.tex_near | 1.44 | 1.31 | 2.50 | - | 1.02 | 2.25 | 1.42 | - | - |
| ICL-NUIM [8] | | | | | | | | | |
| lr_kf0 | 0.35 | 0.37 | 0.37 | - | 0.46 | 1.00 | 0.85 | 0.70 | 0.80 |
| lr_kf1 | 3.96 | 1.04 | 2.40 | 1.60 | 0.72 | 1.16 | 0.64 | 1.10 | 0.98 |
| lr_kf2 | 2.62 | 2.78 | 2.72 | 4.50 | 1.43 | 1.67 | 1.57 | 1.50 | 1.92 |
| lr_kf3 | 1.31 | 1.48 | 2.35 | 4.60 | 1.58 | 1.02 | 1.77 | 1.10 | 1.25 |
| of_kf0 | 2.60 | 4.44 | 6.70 | - | 1.92 | 2.54 | 2.19 | 2.50 | 1.99 |
| of_kf1 | X | X | X | X | 3.15 | 5.64 | 1.89 | 1.30 | 2.25 |
| of_kf2 | 3.17 | 2.18 | 4.54 | 3.10 | 1.59 | 0.97 | 0.87 | 1.50 | 2.20 |
| of_kf3 | 11.1 | 18.03 | 13.50 | 6.50 | 0.91 | 6.94 | 0.96 | 1.30 | 1.84 |

Table 1: **Absolute trajectory error (ATE) RMSE [cm]** (X stands for tracking failure, - stands for not available from the corresponding paper). Each result from ours, ORB-SLAM2 and OpenVSLAM was calculated as the average over 5 executions on each sequence.

ferent planes of similar texture. Notice the network was not trained on the dataset we used to evaluate SLAM, which simulates the practical situation because a trained CNN may not generalize under the different real-world scenarios. This is also the reason why we introduced the graph-cut method in this work which can be considered as a post-processing step to the instance segmentation. Thus, we do not take any output from the CNN as the noise-free measurement of the plane, but further, optimize it within the SLAM framework. We argue that our proposed Sequential Graph-Cut RANSAC (Algo. 1) with plane expanded and re-estimated in SLAM (Sec. 3.3) is an elegant way to reconstruct semantic structure in any monocular SLAM system. Within that, we actually established a procedure where first the potential hypothesis (data points) are separated based on segmentation prior and graph-structured optimization, thereafter the detected geometric primitives are merged based on their geometric attributes during the reconstruction of the scene, on the fly.

4.2 Visual SLAM System with Point and Plane

In this work, we made modification discussed in Sec. 3.3 to OpenVSLAM [23], which is built upon ORB-SLAM2 [18].

Geometric Model Fitting via Energy Minimization. For conducting the energy-based graph-structure optimization, a neighborhood-graph \mathcal{N} has to be built according to the model estimated, where the data point is sampled from. For homography estimation, we define the number of cells along each axis as 8 which is used to divide the image into a grid where the neighborhood-graph is built from. The spatial coherence weight, i.e. the parameter λ in Eq. 2, is set as 0.975. The problem graph construction (line 19 in Algo. 1) within the local optimization (LO) step is used to add energy terms into the function (Eq. 2), where the first energy term (geometric error measure) is added by replacing the standard 0-1 loss with a Gaussian kernel function, which makes the problem close to maximum likelihood estimation, refer to Eq. 1 of [2]. The second energy term (spatial regularization) is added by applying pair-wise energy on a modified Potts model, refer to Eq. 3 of [2]. The error threshold for symmetric transfer error (STE) is set as 2 pixels in this work. The confidence of the inner GC-RANSAC is set as 0.99. For piece-wise planar reconstruction, the neighborhood-graph is constructed via FLANN with the sphere radius set as $r = 2 \cdot \epsilon_d$ (where ϵ_d is the distance threshold discussed in the next paragraph). The confidence of the inner GC-RANSAC is set as 0.99. The spatial coherence weight (λ') is set as 0.6. For all the experiments, the max. iteration number N_{GC} is calculated according to the Eq. 4 in [2], the max. iteration number N of outer RANSAC loop in Algo. 1 is set as 50. **Adaptive Parameter Setting Strategy.** Setting the inlier-outlier threshold is difficult without knowing the scale. Thus, to reduce the dependency on the user-defined threshold, we set two empirical values but adjust it according to the local map scale dynamically: (1) the distance threshold ϵ_d which decides if a point landmark is close

| Thread | Ours | ORB-SLAM2 [18] | OpenVSLAM [23] |
|--|--------------|----------------|----------------|
| Tracking | 16.12 | 20.43 | 19.47 |
| Local Mapping | 83.23 | 110.82 | 105.45 |
| Functionality (refer to Fig. 2) | Ours | | |
| Instance Planar Segmentation | | | 33.11 |
| System (Mono) Initialization | | | 5.37 |
| Point-Plane Map Initialization | | | 22.14 |
| Non-Planar Map Point Cluttering | | | 0.56 |
| New Plane Detection | | | 2.07 |
| Plane Merging/Expanding | | | 0.92 |
| Plane Re-estimation | | | 1.11 |
| Point-Plane Refinement | | | 0.08 |
| Local BA | | | 59.85 |

Table 2: **Runtime analysis [ms] (mean value evaluated on dataset TUM RGB-D [22]: fr3_st.tex_far)** of our planar SLAM system compared to original ORB-SLAM2 and OpenVSLAM, **under monocular setting**, using a desktop PC with an Intel Xeon(R) E-2146G 12 cores CPU @ 3.50GHz, 32GB RAM. The PlaneSegNet is evaluated on a standard GPU of NVIDIA GTX 1080 Ti.

enough to a plane, (2) the residual error ϵ_{Π} which decides if a plane equation fitted is optimal enough (here, ϵ_{Π} is equivalent to the ϵ_L in Algo. 1 which is the residual of a model). Given $\epsilon_d = 0.02$ and $\epsilon_{\Pi} = 0.01$, for monocular mode, the local map scale is estimated as the inverse of median depth of current keyframe during tracking, thereafter both thresholds ϵ_d and ϵ_{Π} are normalized by the local map scale on the fly. For RGB-D mode, the local map scale is calculated as the average value of the sum of the world position of all the point landmarks, this step only needs to be conducted at the beginning of the data sequence. Moreover, the geometric thresholds mentioned in the local mapping thread (Sec. 3.3): $T_{\theta} = 0.8$ (decides if two normals are parallel enough) and $T_d = 10 \cdot \epsilon_d$ (decides if two planes are close enough). We found out above mentioned thresholds give our SLAM system stable performance on various data sequences during evaluation, without much effort on parameter fine-tuning.

Benchmarking. The TUM RGB-D benchmark [22] provides indoor sequences under different texture and structure conditions. Thus, we select different levels of complexity for evaluating our planar SLAM system: single plane scenario (fr1_floor, fr3_nst.tex_near), multiple planes scenario (fr3_st.tex_far, fr3_st.tex_near), and scenario of textureless table but with many objects presented (fr1_xyz, fr1_desk, fr2_xyz, fr2_desk), as listed in Table 1. While our monocular planar SLAM system obtains better results of ATE RMSE on most of the sequence compared to the classic feature-based SLAM systems. Notice that most of the relative work is not superior in terms of ATE RMSE compared to e.g. ORB-SLAM2, such as Structure-SLAM (Mono) and ManhattanSLAM (RGB-D) which reported less complete quantitative results in their paper. Nevertheless, Structure-SLAM, ManhattanSLAM, and SP-SLAM utilize not only plane features but also lines, predicted normal maps or MW assumption in their system. Ours, however, presented as a pure point-plane SLAM system without WM assumption whose algorithm is suitable for any feature-based SLAM of monocular or RGB-D, which shows integrating MW assumption or more high-level features (e.g. line) benefit semantic mapping or tracking in the low-texture scenes, but not necessarily improve the accuracy of camera localization, possibly due to different optimization strategies. More important, any output from CNN should not be considered as noise-free input for the semantic SLAM system, as we especially tackled in this work with graph-cut optimization within SLAM workflow. The qualitative results are illustrated in Fig. 4.

The ICL-NUIM RGB-D benchmark [8] is a synthetic indoor dataset that shows a low-contrast and low-texture environment. Thus we lower the FAST threshold for detecting ORB features to 2, which gives the best performance under monocular setting, while RGB-D SLAM is evaluated under the default setting. However, even with a lower FAST threshold, this dataset is difficult for monocular SLAM and our monocular system only works stably on a few of the sequences, and we are not able to initialize the system when

processing more than half of the images of sequence of_kt1. Without establishing 3D plane-plane registration and pose optimization [24] as done in mainstream approaches of RGB-D planar SLAM, our RGB-D system also obtains comparable ATE RMSE accuracy compare to other RGB-D SLAM methods, as shown in Table 1.

Failure Cases and Limitation. Our monocular planar SLAM system depends strongly on the point features, which brings the limitation that no reliable plane can be fitted when there are not enough point landmarks. Our map refiner strategy only keeps high-quality planes, which also omits small planes from the map. Textureless planar scenes such as fr1_desk and fr2_desk, result in fitted planes from point cloud actually associated with objects like books, keyboards, and cups, which undermines the camera localization. The ATE RMSE evaluation of our SLAM systems also strongly depends on the performance of OpenVSLAM. The tracking failure of sequence of_kt1 (monocular) results from a fast-moving and fast-rotating camera and textureless scene. RGB-D SLAM is usually more robust but more 3D point and plane landmarks could result in a large-scale non-linear optimization even in a small scene.

Run-time. The run-time analysis of our monocular SLAM system is reported in Table 2, which shows a more efficient computation time compared to the classic feature-based ORB-SLAM2 and OpenVSLAM. Notice the functionality presented in Table 2 is corresponding to the building block illustrated in Fig. 2, which is embedded as a multi-thread system. Most of the computation is used for initialization and local BA. After the system is initialized, the instance planar segmentation only needs to be conducted on every inserted keyframe. It is very fast to detect new planes using our proposed sequential multi-plane fitting method (avg. 2 ms) when a new keyframe is inserted into the map, and it is very fast to merge, expand, and refine the existing plane structures (avg. 1 ms).

5 CONCLUSION

Our work presented a robust building block of a feature-based SLAM framework with an extended plane detector, with special care in taking instance plane segmentation as noisy "sensor" input and further optimizing it during geometric primitives reconstruction. With the dynamically adjusted thresholds, our proposed multi-plane reconstructor can be applied to various indoor scenarios without much effort in parameter fine-tuning. Comprehensive quantitative results are reported in this work. Future works can explore other types of features such as line segment and vanish point, or utilize planar SLAM in the urban environment and driving scenario. The joint pose-graph optimization of different geometric primitives is of interest.

REFERENCES

- [1] D. Barath, M. Ivaschkin, and J. Matas. Progressive napsac: sampling from gradually growing neighborhoods. *arXiv preprint arXiv:1906.02295*, 2019.
- [2] D. Barath and J. Matas. Graph-cut ransac. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [3] D. Barath and J. Matas. Progressive-x: Efficient, anytime, multi-model fitting algorithm. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3780–3788, 2019.
- [4] Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Transactions on pattern analysis and machine intelligence*, 23(11):1222–1239, 2001.
- [5] O. Chum, J. Matas, and J. Kittler. Locally optimized ransac. In *Joint Pattern Recognition Symposium*, pp. 236–243. Springer, 2003.
- [6] J. Engel, T. Schöps, and D. Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pp. 834–849. Springer, 2014.
- [7] M. A. Fischler and R. C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- [8] A. Handa, T. Whelan, J. McDonald, and A. J. Davison. A benchmark for rgb-d visual odometry, 3d reconstruction and slam. In *2014 IEEE international conference on Robotics and automation (ICRA)*, pp. 1524–1531. IEEE, 2014.
- [9] M. Hsiao, E. Westman, G. Zhang, and M. Kaess. Keyframe-based dense planar slam. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5110–5117. IEEE, 2017.
- [10] H. Isack and Y. Boykov. Energy-based geometric multi-model fitting. *International journal of computer vision*, 97(2):123–147, 2012.
- [11] M. Kaess. Simultaneous localization and mapping with infinite planes. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4605–4611. IEEE, 2015.
- [12] G. H. Lee, F. Fraundorfer, and M. Pollefeys. Mav visual slam with plane constraint. In *2011 IEEE International Conference on Robotics and Automation*, pp. 3139–3144. IEEE, 2011.
- [13] H. Li. Two-view motion segmentation from linear programming relaxation. In *2007 IEEE conference on computer vision and pattern recognition*, pp. 1–8. IEEE, 2007.
- [14] Y. Li, N. Brasch, Y. Wang, N. Navab, and F. Tombari. Structure-slam: Low-drift monocular slam in indoor environments. *IEEE Robotics and Automation Letters*, 5(4):6583–6590, 2020.
- [15] C. Liu, K. Kim, J. Gu, Y. Furukawa, and J. Kautz. PlanerCNN: 3d plane detection and reconstruction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4450–4459, 2019.
- [16] C. Liu, J. Yang, D. Ceylan, E. Yumer, and Y. Furukawa. Planenet: Piecewise planar reconstruction from a single rgb image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2579–2588, 2018.
- [17] M. Muja and D. G. Lowe. Fast approximate nearest neighbors with automatic algorithm configuration. *VISAPP (1)*, 2(331-340):2, 2009.
- [18] 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*, 33(5):1255–1262, 2017.
- [19] J. Rambach, P. Lesur, A. Pagani, and D. Stricker. Slamcraft: Dense planar rgb monocular slam. In *2019 16th International Conference on Machine Vision Applications (MVA)*, pp. 1–6. IEEE, 2019.
- [20] C. Raposo and J. P. Barreto. π match: Monocular vslam and piecewise planar reconstruction using fast plane correspondences. In *European Conference on Computer Vision*, pp. 380–395. Springer, 2016.
- [21] R. F. Salas-Moreno, B. Glocken, P. H. Kelly, and A. J. Davison. Dense planar slam. In *2014 IEEE international symposium on mixed and augmented reality (ISMAR)*, pp. 157–164. IEEE, 2014.
- [22] 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.
- [23] S. Sumikura, M. Shibuya, and K. Sakurada. Openvslam: A versatile visual slam framework. In *Proceedings of the 27th ACM International Conference on Multimedia*, pp. 2292–2295, 2019.
- [24] Y. Taguchi, Y.-D. Jian, S. Ramalingam, and C. Feng. Point-plane slam for hand-held 3d sensors. In *2013 IEEE international conference on robotics and automation*, pp. 5182–5189. IEEE, 2013.
- [25] X. Wang, M. Christie, and E. Marchand. Relative pose estimation and planar reconstruction via superpixel-driven multiple homographies. In *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems, IROS'20*, 2020.
- [26] Y. Xie, J. Rambach, F. Shu, and D. Stricker. Planesegnet: Fast and robust plane estimation using a single-stage instance segmentation CNN. *arXiv preprint arXiv:2103.15428*, 2021.
- [27] S. Yang and S. Scherer. Monocular object and plane slam in structured environments. *IEEE Robotics and Automation Letters*, 4(4), 2019.
- [28] R. Yunus, Y. Li, and F. Tombari. Manhattanslam: Robust planar tracking and mapping leveraging mixture of manhattan frames. *arXiv preprint arXiv:2103.15068*, 2021.
- [29] X. Zhang, W. Wang, X. Qi, Z. Liao, and R. Wei. Point-plane slam using supposed planes for indoor environments. *Sensors*, 19(17):3795, 2019.
- [30] M. Zuliani, C. S. Kenney, and B. Manjunath. The multiransac algorithm and its application to detect planar homographies. In *IEEE International Conference on Image Processing 2005*. IEEE, 2005.