

"Vision based SLAM in dynamic scenes"

2012-Dec-20

Name of Principal Investigators (PI and Co-PIs):

- e-mail address : eletp@nus.edu.sg
- Institution : National University of Singapore
- Mailing Address : Department of Electrical and Computer Engineering
- Phone : +65-6516-2130
- Fax : +65-6779-1103

Period of Performance: 09/26/2011 – 09/25/2012

Abstract: In this project, we studied the vision-based simultaneous localization and mapping (SLAM) problem from a novel perspective. We employed multiple independently moving cameras, while conventional studies are limited with a single camera (or a multi-camera rig where the relative positions between cameras are fixed). Our flexible configuration of cameras makes this algorithm applicable to robot teams, which also makes this study the world's first vision based SLAM algorithm for robot teams. Furthermore, the collaboration among multiple cameras allows us to deal with challenging dynamic scenes which make most of previous algorithms fail. This work was accepted for publication at the IEEE Transaction of Pattern Analysis and Machine Intelligence (TPAMI), with impact factor 5.9, #1 in all engineering and AI.

Introduction: Please refer to the attached technique paper.

Experiment: Please refer to the attached technique paper.

Results and Discussion: Please refer to the attached technique paper.

List of Publications and Significant Collaborations that resulted from your AOARD supported project:

Danping Zou, Ping Tan, CoSLAM: Collaborative Visual SLAM in Dynamic Environments, IEEE Transaction on Pattern Analysis and Machine Intelligence, Accepted for publication.

Attachments:

Report Documentation Page			Form Approved OMB No. 0704-0188		
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE 17 JAN 2013		2. REPORT TYPE Final		3. DATES COVERED 26-09-2011 to 25-09-2012	
4. TITLE AND SUBTITLE Vision based SLAM in dynamic scenes			5a. CONTRACT NUMBER FA23861114118		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S) Ping Tan			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) National University of Singapore,4 Engineering Drive 3,Singapore 117576,Singapore,SG,117576			8. PERFORMING ORGANIZATION REPORT NUMBER N/A		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) AOARD, UNIT 45002, APO, AP, 96338-5002			10. SPONSOR/MONITOR'S ACRONYM(S) AOARD		
			11. SPONSOR/MONITOR'S REPORT NUMBER(S) AOARD-114118		
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT In this project, we studied the vision-based simultaneous localization and mapping (SLAM) problem from a novel perspective. We employed multiple independently moving cameras, while conventional studies are limited with a single camera (or a multi-camera rig where the relative positions between cameras are fixed). Our flexible configuration of cameras makes this algorithm applicable to robot teams, which also makes this study the world's first vision based SLAM algorithm for robot teams. Furthermore, the collaboration among multiple cameras allows us to deal with challenging dynamic scenes which make most of previous algorithms fail. This work was accepted for publication at the IEEE Transaction of Pattern Analysis and Machine Intelligence (TPAMI), with impact factor 5.9, #1 in all engineering and AI.					
15. SUBJECT TERMS Guidance, Navigation,Guidance, and Control					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT Same as Report (SAR)	18. NUMBER OF PAGES 16	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

CoSLAM: Collaborative Visual SLAM in Dynamic Environments

Danping Zou and Ping Tan, *Member, IEEE*,

Abstract—This paper studies the problem of vision-based simultaneous localization and mapping (SLAM) in dynamic environments with multiple cameras. These cameras move independently and can be mounted on different platforms. All cameras work together to build a global map, including 3D positions of static background points and trajectories of moving foreground points. We introduce inter-camera pose estimation and inter-camera mapping to deal with dynamic objects in the localization and mapping process. To further enhance the system robustness, we maintain the position uncertainty of each map point. To facilitate inter-camera operations, we cluster cameras into groups according to their view overlap, and manage the split and merge of camera groups in real-time. Experimental results demonstrate that our system can work robustly in highly dynamic environments and produce more accurate results in static environments.

Index Terms—Visual SLAM, Swarm, Dynamic Environments, Structure-from-Motion

1 INTRODUCTION

Many vision based SLAM (simultaneous localization and mapping) systems [23], [10], [18] have been developed, and have shown their remarkable performance on mapping and localization in real-time. Recent works took further steps to provide high level scene understanding [20], or to improve the system accuracy and robustness, such as ‘loop closure’ [16], ‘re-localization’ [36], and dense depth map reconstruction [22]. These works have made cameras become more and more favorable sensors for SLAM systems.

Existing vision-based SLAM systems mainly focus on navigation in static environments with a single camera. However, the real world is full of moving objects. Although there are robust methods to detect and discard dynamic points by treating them as outliers [18], [5], conventional SLAM algorithms tend to fail when the portion of moving points is large. Further, in dynamic environments, it is often important to reconstruct the 3D trajectories of the moving objects [38], [35] for tasks such as collision detection and path planning. This 3D reconstruction of dynamic points can hardly be achieved by a single camera.

To address these problems, we present a collaborative visual SLAM system using multiple cameras. The relative positions and orientations between cameras are allowed to change over time, which means cameras can be mounted on different platforms that move independently. This setting is different from existing SLAM systems with a stereo camera [23], [25] or a multi-camera rig [17] where all cameras are fixed on a single platform. Our camera configuration makes the system applicable to the following interesting cases:

1) wearable augmented reality [5], where multiple cameras are mounted on different parts of the body; 2) robot teams [4], [34], [1], where multiple robots work in the same environment and each carries a single camera because of limited weight and energy capacity, e.g. micro air vehicles (MAVs) [40].

In our system, we use images from different cameras to build a global 3D map. We maintain the position uncertainty of each map point by a covariance matrix, and iteratively refine the map point position whenever a new observation is available. This design increases the system robustness and accuracy in dealing with complex scenes. Further, we classify map points as dynamic or static at every frame by analyzing their triangulation consistency. False points caused by incorrect feature matching are also detected and removed. For robust localization in dynamic environments, we use both dynamic and static points to simultaneously estimate the poses of all cameras with view overlap. We divide cameras into groups according to their view overlap. Cameras within the same group share a common view, and work collaboratively for robust mapping and localization. Groups can split or merge when cameras separate or meet.

Our system was tested in static and dynamic, indoor and outdoor scenes. Experimental results show that our method is more accurate and robust than existing single camera based SLAM methods. Our system succeeds in highly dynamic environments, and is able to reconstruct the 3D trajectories of moving objects. The system currently runs at approximately 38 ms per frame for three cameras. In the next section we shall briefly review existing visual SLAM methods and discuss our work in detail.

• D. Zou and P. Tan are with the National University of Singapore.
E-mail: dannis.zou@gmail.com; eletp@nus.edu.sg.

2 RELATED WORK

Visual SLAM with a single camera There are mainly two types of methods for single camera based visual SLAM. One is based on the structure-from-motion (SFM) technique [15]. Royer et al. [26] first reconstructed a 3D map of the scene offline from a learning sequence, and then estimated the camera pose in real-time by referring to that map. Mouragnon et al. [21] proposed a local bundle adjustment method so that mapping and pose update can run in nearly real-time. Klein et al. [18] put the time-consuming bundle adjustment into an asynchronous thread, and made the system much faster.

The second type of methods model SLAM as a Bayesian inference problem, and solve it through the Extended Kalman Filter [9], [10]. In [30], line features were used to complement point features to improve the matching robustness. Eade et al. [12] applied particle filter and a top-down approach to handle a relatively large number of landmarks (hundreds) more efficiently. To improve the robustness of SLAM, Williams et al. [36] proposed a re-localization method to recover the SLAM system from tracking failures.

Strasdat et al. [31], [33] compared both types of methods, and concluded that the SFM based methods produce more accurate results per unit computing time, while filter based methods could be more efficient when processing resource is limited.

These methods often do not consider dynamic scenes. Some of them, such as [18], [5], detected and discarded dynamic points as outliers. However, this approach tends to fail when the portion of dynamic points is large. Some more recent methods, such as [24], [19], applied multi-body SFM to deal with dynamic environments. However, this approach is only applicable to rigid moving objects, and the 3D reconstruction of moving points is up to a scaling ambiguity [24]. In comparison, we propose to solve the SLAM problem in dynamic scenes with multiple independently moving cameras. Our approach can reconstruct the full 3D of dynamic points within the scene map.

Visual SLAM with multiple cameras Nister et al. [23] proposed a visual odometry system with a stereo rig. Their system was much like a SFM-based single camera SLAM system with an additional camera to generate map points at every frame. They also pointed out the narrow base line between stereo cameras could affect the map quality. To address this problem, Paz et al. [25] separated close and far 3D points, and used far points to estimate camera rotation only. To obtain wider FOV, Kaess et al. [17] mounted multiple cameras in a rig facing different directions to combine the advantages of omnidirection vision [37] and monocular vision. Castle et al. [5] used multiple cameras distributed freely in a static environment, where each camera was processed by an independent

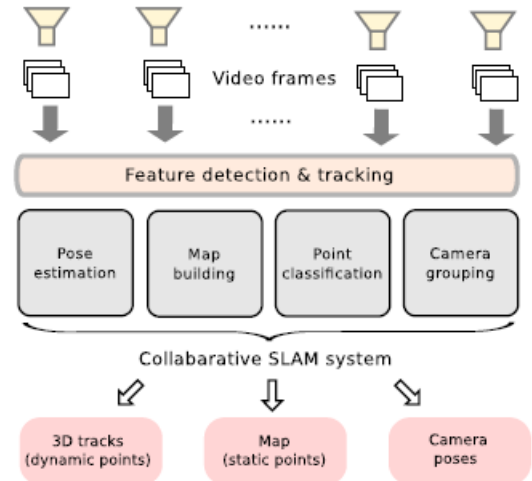


Fig. 1. CoSLAM system architecture.

single camera based SLAM system. A camera could be localized according to different maps by registering its feature points to the other map points.

These methods still focus on static scenes and do not take full advantage of multiple cameras. Further, the relative positions of their cameras are often fixed. In comparison, we allow cameras to move independently and achieve more flexible system design. For example, our cameras can be mounted on different robots for robot team applications.

SLAM in dynamic environments Existing works on SLAM in dynamic environments mostly use filter-based methods and have been successfully applied to SLAM problems with sensors such as laser scanners [14], [38], [35] and radar systems [2], [3]. In this work, we aim to use cameras to address the SLAM problem. Compared with other sensors, cameras are passive, compact and energy efficient, which have important advantages for micro robots with limited weight and energy capacity (such as MAVs [40], [1]). To the best of our knowledge, this work is the first visual SLAM solution in a dynamic environment with multiple cameras moving independently. This method could be applied to emerging swarm robotics applications [1], [27].

3 SYSTEM OVERVIEW

The intrinsic parameters of all our cameras are calibrated beforehand. Our collaborative SLAM system treats each camera as a sensor input, and incorporates all inputs to build a global map, and simultaneously computes the poses of all cameras over time. The system framework is illustrated in Figure 1. The system detects and tracks feature points at every frame, and feeds them to the four SLAM components. We use Kanade-Lucas-Tomasi(KLT)[28] tracker for both feature detection and tracking because of its good balance between efficiency and robustness. However, there is no restriction to use other feature detectors and trackers such as the ‘active matching’ [6].

The four SLAM components are ‘camera pose estimation’, ‘map building’, ‘point classification’, and ‘camera grouping’. The main pipeline of our system follows conventional sequential structure-from-motion (SfM) methods. We assume that all cameras look at the same initial scene to initialize the system. After that, the ‘camera pose estimation’ component computes camera poses at every frame by registering the 3D map points to 2D image features. From time to time, new map points are generated by the ‘map building’ component. At every frame, points are classified into different types by the ‘point classification’ component. The system maintains the view overlap information among cameras throughout time. The ‘camera grouping’ component separates cameras into different groups, where cameras with view overlap are in the same group. These groups could merge and split when cameras meet or separate. In the following section, we shall describe these four components in detail.

4 CAMERA POSE ESTIMATION

Our system alternatively uses two different methods for camera pose estimation: intra-camera pose estimation and inter-camera pose estimation. In the former, each camera works independently, where tracked feature points from a camera are registered with static map points to compute its pose. In dynamic environments, the number of static map points could be small, or the static points are distributed within a small image region, which can make the intra-camera pose estimation fail. In such a case, we switch to the inter-camera pose estimation method that uses both static and dynamic points to simultaneously obtain poses for all cameras.

4.1 Intra-camera Pose Estimation

If the camera intrinsic parameters are known, the camera pose $\Theta = (\mathbf{R}, \mathbf{t})$ can be computed by minimizing the reprojection error (the distance between the image projection of 3D map points and their corresponding image feature points), namely,

$$\Theta^* = \arg \min_{\Theta} \sum_i \rho(\|\mathbf{m}_i - \mathcal{P}(\mathbf{M}_i, \Theta)\|). \quad (1)$$

where $\mathcal{P}(\mathbf{M}_i, \Theta)$ is the image projection of the 3D point \mathbf{M}_i , \mathbf{m}_i is the image feature point registered to \mathbf{M}_i , $\|\cdot\|$ measures the distance between two image points. i is an index of feature points. The M-estimator $\rho: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is the Tukey bi-weight function [39] defined as

$$\rho(x) = \begin{cases} t^2/6(1 - [1 - (\frac{x}{t})^2]^3) & \text{if } |x| \leq t \\ t^2/6 & \text{otherwise.} \end{cases} \quad (2)$$

Assuming that the error of feature detection and tracking obeys a Gaussian distribution $\mathcal{N}(0, \sigma^2 \mathbf{I})$, we set the threshold t in $\rho(\cdot)$ as 3σ . Equation (1) is

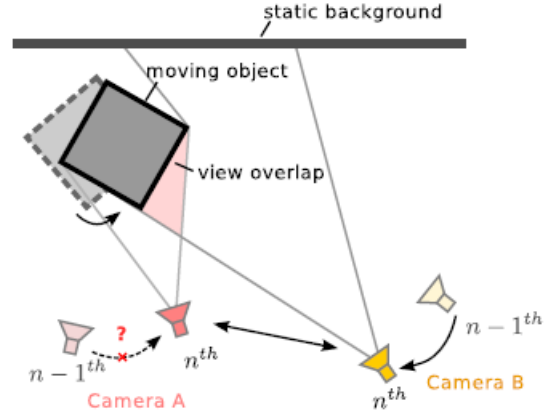


Fig. 2. The pose of camera A at the n^{th} frame cannot be estimated from its previous pose, since it observes only the moving object (the grey square). However, its relative pose with respect to camera B can be determined. So the absolute pose of camera A can be computed.

minimized by the iteratively re-weighted least squares (IRLS) method, where Θ is initialized according to the camera pose at the previous frame. At each iteration of the IRLS, the Levenberg-Marquart algorithm is used to solve the non-linear least square problem, where Θ is parameterized in Lie algebra $se(3)$ as [32].

4.2 Inter-camera Pose Estimation

When the number of visible static points is small, or the static points are located in a small image region, the intra-camera pose estimation is unstable and sometimes fails. Fortunately, points on moving objects give information about the relative camera poses. Figure 2 provides such an illustration. The pose of camera A at the n^{th} frame cannot be decided on its own, since it only observes the moving object (the grey square). However, its relative pose with respect to the camera B can be decided. We can therefore use both static and dynamic points together to decide all camera poses simultaneously.

Actually, the 3D coordinates of dynamic points can only be computed when the camera poses are already known. Hence, our system in fact simultaneously estimates both the camera poses and the 3D positions of dynamic points. We formulate the inter-camera pose estimation problem as an minimization of reprojection error,

$$\{\Theta_c\}^* = \arg \min_{\mathbf{M}_D, \{\Theta_c\}} \sum_c \left\{ \sum_{i \in S} v_i^c \rho(\|\mathbf{m}_i - \mathcal{P}(\mathbf{M}_i, \Theta_c)\|) + \sum_{j \in D} v_j^c \rho(\|\mathbf{m}_j - \mathcal{P}(\mathbf{M}_j, \Theta_c)\|) \right\}. \quad (3)$$

Here, c is an index of cameras, S and D are the set of ‘static’ and ‘dynamic’ map points. v_i^c represents the

visibility of the i -th map point at camera c (1 for visible, 0 otherwise).

The difference between the ‘intra-camera pose estimation’ and the ‘inter-camera pose estimation’ lies in the second term of Equation (3), where the dynamic points are included in the objective function. The relative poses between cameras are therefore enforced by minimizing the reprojection error of the dynamic points. Hence, our system can determine the poses of cameras where few static points are visible. As the cameras need to have view overlap, we only apply inter-camera pose estimation to cameras within the same group. We refer the reader to Section 6 for more details about camera grouping.

The optimization of Equation (3) is also solved by the IRLS. The camera poses and 3D positions of dynamic points are initialized from the previous frame. We call intra-camera pose estimation by default. We call the inter-camera pose estimation only when the number of dynamic points are greater than that of static points, or the area covered by the convex hull of static feature points is less than 20% of the image area.

5 MAP MAINTENANCE

Unlike previous SFM based SLAM systems [18], [5], where only the 3D position is recovered for each map point, we further maintain the position uncertainty of each map point to help point registration and distinguishing static and dynamic points. Specifically, we recover a probability distribution of the map point position, which is represented by a Gaussian function $\mathcal{N}(\mathbf{M}_i, \Sigma_i)$. \mathbf{M}_i is the triangulated position. The covariance matrix $\Sigma_i \in \mathbb{R}^{3 \times 3}$ measures the position uncertainty.

To facilitate computation, for each feature point, we keep a pointer directing to its corresponding map point. Similarly, for each map point, we also keep pointers directing to its corresponding image feature points in each view, and store the local image patches centered at these feature points. We downsize the original input image to its 30% size, and take a patch of 11×11 pixels. We further keep track of the frame numbers when a map point is generated or becomes invisible.

5.1 Position Uncertainty of Map Points

When measuring the uncertainty in map point positions, we only consider the uncertainty in feature detection and triangulation. In principle, we could also include the uncertainty in camera positions. We assume the feature detection error follows Gaussian distribution $\mathcal{N}(0, \sigma^2 \mathbf{I})$. The position uncertainty of a 3D map point is described by the covariance computed as

$$\Sigma_i = (\mathbf{J}_i^T \mathbf{J}_i)^{-1} \sigma^2, \quad (4)$$

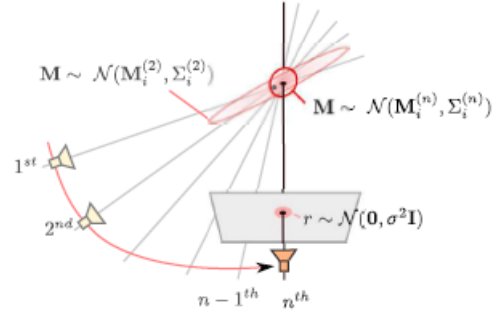


Fig. 3. Both the mean position and position covariance matrix are updated when a new observation is coming.

where $\mathbf{J}_i \in \mathbb{R}^{(2k) \times 3}$ is the Jacobian of the camera projection function that maps a 3D map point to its 2D image coordinates in all views, and k denotes the number of views used for triangulation.

When there is a new image observation $\mathbf{m}_i^{(n+1)}$ of a map point, we can quickly update its 3D position with Kalman Gain

$$\mathbf{M}_i^{(n+1)} = \mathbf{M}_i^{(n)} + \mathbf{K}_i [\mathcal{P}_{n+1}(\mathbf{M}_i^{(n)}) - \mathbf{m}_i^{(n+1)}]. \quad (5)$$

Here, $\mathcal{P}_{n+1}(\mathbf{M}_i^{(n)})$ computes the image projection of $\mathbf{M}_i^{(n)}$ in the $(n+1)$ th frame. The Kalman Gain $\mathbf{K}_i \in \mathbb{R}^{3 \times 2}$ is computed as

$$\mathbf{K}_i = \Sigma_i^{(n)} \hat{\mathbf{J}}_i^T \left(\sigma^2 \mathbf{I} + \hat{\mathbf{J}}_i^T \Sigma_i^{(n)} \hat{\mathbf{J}}_i \right)^{-1}. \quad (6)$$

Here, $\hat{\mathbf{J}}_i \in \mathbb{R}^{2 \times 3}$ is the Jacobian of $\mathcal{P}_{n+1}(\cdot)$ evaluated at $\mathbf{M}_i^{(n)}$. The triangulation uncertainty can meanwhile be updated by

$$\Sigma_i^{(n+1)} = \Sigma_i^{(n)} - \mathbf{K}_i \hat{\mathbf{J}}_i \Sigma_i^{(n)}. \quad (7)$$

We illustrate the idea of this refinement in Figure 3.

Unlike SLAM algorithms based on Kalman filter, such as [9], [11], which spend a high computational cost $O(N^2)$ (N is the number of map points) to maintain the covariance matrix for both camera states and the whole map states (including correlation between positions of different points), we only maintain the triangulation uncertainty for each individual map point. The computation cost of our method is only $O(N)$. Further, the computation is independent at each point, which enables parallel computation to achieve better efficiency. Maintaining position uncertainty is important for the following map operations such as point classification and registration described in Section 5.3 and Section 5.4. For example, the point classification benefits from the position uncertainty. Even for static points, their reconstructed positions are always changing over time due to triangulation uncertainties. (This can be seen clearly in our supplementary videos.) With the position covariance matrix Σ_i , we can distinguish static and dynamic points better. Though we could also maintain the covariance among different points and uncertainty of camera positions, we do not include them for efficiency consideration.

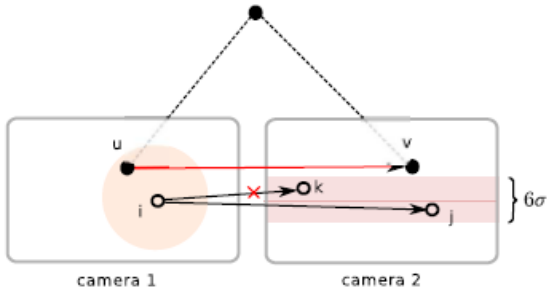


Fig. 4. Guided feature matching between two cameras. u, v are the projections of a map point. We will prefer the match $i \leftrightarrow j$ than $i \leftrightarrow k$ for its better consistency to $u \leftrightarrow v$.

5.2 Map Points Generation

We propose two methods to generate new map points. The first one, ‘intra-camera mapping’, reconstructs static map points from feature tracks in each individual camera. To deal with moving objects, we propose the second method, ‘inter-camera mapping’, to generate map points from corresponding points across cameras within the same group (i.e. cameras with view overlap).

5.2.1 Intra-camera mapping

Previous methods usually use key frames to generate new map points [18], [26], where feature points in selected key frames are matched and triangulated. Since all feature points are tracked over time in our system, it is not necessary to select key frames to match those feature points again. If there are unmapped feature tracks (whose pointers to map points are NULL) that are long enough ($> N_{min}$ frames), we use the beginning and the end frames of this track to triangulate a 3D point. Once the 3D position is computed, the covariance can also be evaluated by Equation (4). We then check the reprojection error at all frames of the feature track. If the Mahalanobis distance (described in Section 5.3) between the projection and the feature point is smaller than θ for all frames, a new map point is generated and marked as ‘static’.

5.2.2 Inter-camera mapping

Inter-camera mapping is applied to unmapped feature points only. We match the image features between different cameras by zero-mean normalized cross correlation (ZNCC). To avoid ambiguous matches, the corresponding points are searched only in a narrow band within 3σ distance to the epipolar line. Only matches with $\text{ZNCC} > \tau_{ncc}$ are considered. We further use the correspondence of existing map points as seeds to guide matching - a pair of feature points is not considered to be matched, if it has a very difference disparity vector from that of the nearest seed.

As shown in Figure 4, suppose we want to find the corresponding point for the unmapped feature point

i . There are two candidates k and j within the band of the epipolar line. u is the closest mapped feature point to i . If it is within ϕ_r pixels distance to i , we use the disparity vector D_{uv} to guide feature matching. We compare the difference between the disparity vectors. A candidate with very different disparity from the seed is discarded, e.g. candidate k is removed in Figure 4 because $\|D_{uv} - D_{ik}\| > \phi_d$. The best match is then obtained from the remaining candidates by the winner-take-all strategy.

After matching feature points between cameras, we triangulate the corresponding points to generate new map points. Exhaustive feature matching in all possible camera pairs is inefficient. We construct a graph for cameras within the same group where cameras are linked according to their view overlap. We select a spanning tree of the graph and only match features between cameras that are connected by the spanning tree edges. More details of this spanning tree are discussed in Section 6. Inter-camera mapping are called every 5 frame in our system when dynamic points are detected.

5.3 Point Registration

At every frame, we need to associate each map point with its incoming observations - newly detected image feature points from different cameras. Many of the feature points are registered to a map point through feature tracking. We further process the remaining unmapped feature points. For these points, we only consider *active* map points which are static and have corresponding feature points within the most recent N_{rec} frames. These map points are cached in our system for fast access. For each unmapped feature point detected at the current frame, we go through these *active* map points for registration.

We project *active* map points to the images, and compare the image patches centered at the projections with that of the feature point through ZNCC. Considering the uncertainty in map point position and feature detection, the projected position of a map point M_i should satisfy a Gaussian distribution $\mathcal{N}(m_i, c_i)$, where the covariance $c_i = \hat{J}_i \Sigma_i \hat{J}_i^T + \sigma I$, $\hat{J}_i \in \mathbb{R}^{2 \times 3}$ is the Jacobian of the image projection $\mathcal{P}(\cdot)$ evaluated at M_i . The Mahalanobis distance is computed as

$$D^2(m_j, m_i) = (m_j - m_i)^T c_i^{-1} (m_j - m_i). \quad (8)$$

We only consider the feature point m_j which has the smallest Mahalanobis distance to m_i . We then check the ZNCC score between M_i and m_j . To alleviate problems caused by perspective distortion, when selecting the image patch for M_i , we choose the one stored from the nearest camera. m_j is discarded if its ZNCC score $< \tau_{ncc}$. We further traverse back along the feature track of m_j to check if its previous positions are also nearby to the projections of M_i . If

the Mahalanobis distances between them in all frames are smaller than θ , then the m_j is registered to M_i .

Once an unmapped feature point is registered to a 3D map point, the 3D position and the position covariance of this map point can be updated based on the new observation. However, the new observation obtained from point registration, unlike those obtained from feature tracking, usually deviates largely from the previous observations (e.g. this new observation often comes from a different camera). It will lead to inaccurate estimation if we use the iterative refinement described in Equation (5) and Equation (7). In such case, we retriangulate the 3D position of this map point with all observations and recompute the covariance by Equation (4). To reduce the computational cost, we select only two observations from the feature track in each camera for retriangulation, which have the largest viewpoint changes.

5.4 Point Classification

At every frame, we need to distinguish ‘dynamic’ points on moving objects and ‘static’ points on the background. A naïve method for this classification is to threshold the position variation of the 3D map points. It is, however, difficult to set a common threshold for different scenes. Further, in our system, the positions of static points are also changing over time, since their positions are updated whenever new observations are available. Especially, for static points in the distant background, their positions may change significantly as cameras move from far to near.

We instead distinguish static and dynamic points by the reprojection error, which can be easily measured on the image plane. For a dynamic point, if we project its 3D position at the $(n-1)^{\text{th}}$ frame to the n^{th} frame, since the 3D position actually changes, the projection should be distant from its tracked corresponding feature points. In other words, if we use image feature points from different (or the same) time to triangulate the 3D position of a dynamic point, the reprojection error should be large (or small). In comparison, if the point is static, the reprojection error should be always small, no matter if we use image feature points from the same or different time. Based on this observation, we design a process to distinguish ‘static’ and ‘dynamic’ points. The whole process is illustrated in Figure 5. We use ‘false’ points to denote map points generated from incorrect correspondence. An intermediate state ‘uncertain’ is also introduced for points need further investigation.

Initially, we consider all points as static. At every frame, we check the reprojection errors of all ‘static’ points. The projected position of a static map point obeys Gaussian distribution $\mathcal{N}(m_i, c_i)$. The Mahalanobis distance between corresponding feature points and m_i should be less than θ . If the tracked feature point has larger Mahalanobis distance, the

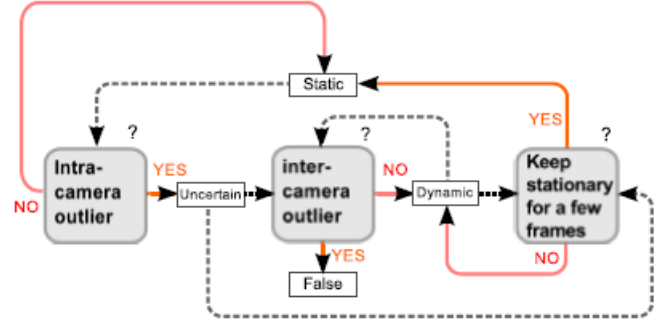


Fig. 5. Map point classification. The pipeline for classifying a point into four types : ‘static’, ‘dynamic’, ‘false’ or ‘uncertain’.

map point is likely to be ‘dynamic’ or ‘false’. We consider these points are *intra-camera outliers*, i.e. outliers for intra-camera triangulation. We mark them as ‘uncertain’ for the next step of classification.

An uncertain point could be either ‘dynamic’ or ‘false’. To distinguish them, we re-triangulate its 3D position M_i with its tracked feature points in the same frame from different cameras. If the Mahalanobis distances of all these feature points to the projection of M_i are smaller than θ , we consider the map point as ‘dynamic’. Otherwise, it is an *inter-camera outlier*, i.e. an outlier for inter-camera triangulation. We consider *inter-camera outliers* as ‘false’ points caused by incorrect feature matching. Note that the 3D positions of dynamic points are naturally updated over time during point classification. Hence, our system is able to produce the 3D trajectories of moving points.

A dynamic point may become static if the object stops moving. Hence, we project the current 3D position of a dynamic point to the previous frames. If the projection is close to its tracked feature points (Mahalanobis distance $< \theta$) for N_{min} of continuous frames, we consider this point as ‘static’.

Figure 6 shows two video frames of a camera from the ‘sitting man’ example, where all the points on the body should be dynamic. We use green and blue points to visualize static and dynamic points. Though all map points were marked as ‘static’ initially, during the SLAM process, our map point classification component successfully differentiated dynamic and static points. This can be seen from the right in Figure 6, where all points on the body were marked in blue.

6 CAMERA GROUPING

The inter-camera operations, e.g. mapping and pose estimation, can be only applied to cameras with view overlap. As cameras move independently, the view overlap among cameras changes over time. In this section, we describe our method to identify and manage camera groups, where cameras with view overlap are in the same group.

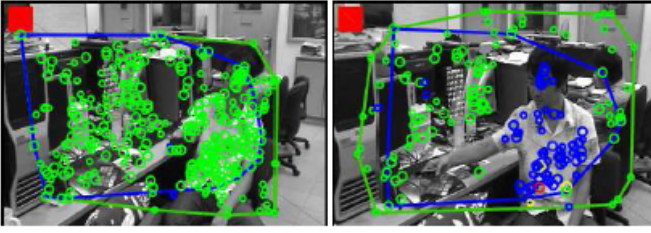


Fig. 6. Map point classification. Green and blue points indicate static and dynamic map points. Left: during the initialization, all points were marked as ‘static’. Right: our system correctly identified the dynamic points by map classification.

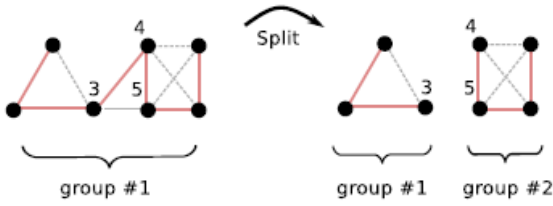


Fig. 7. Camera grouping and splitting. Each node is a camera and each edge links two cameras with view overlap. The solid edges are those on the spanning trees.

6.1 Grouping and Splitting

Since we store a pointer to the corresponding 3D point for each mapped feature point, we can quickly count the number of common map points N_{ij} between two cameras i, j . We construct an undirected graph where the nodes represent the cameras. If $N_{ij} > 0$, we connect the camera i and j by an edge weighted by N_{ij} . A connected component in this graph forms a camera group. The inter-camera operations are only applied to cameras in the same group.

As discussed in Section 5.2.2, to improve the efficiency of inter-camera mapping, we do not match feature points between all camera pairs with view overlap. Instead, we extract a spanning tree for each camera group with maximum weight, and only match feature points between cameras if the edge connecting them is on the selected spanning tree.

A camera group will split when any camera in it moves away and does not have view overlap with others. Such a case is illustrated in Figure 7. The edges between camera 3, 4 and 3, 5 are removed as they have no common feature points any more. The original graph is therefore separated into two connected components, each of which forms a new camera group at the current frame. A real example is provided in Figure 8. Four cameras were split into two groups, where the red and green cameras share a common view, and the blue and yellow cameras share another common view.

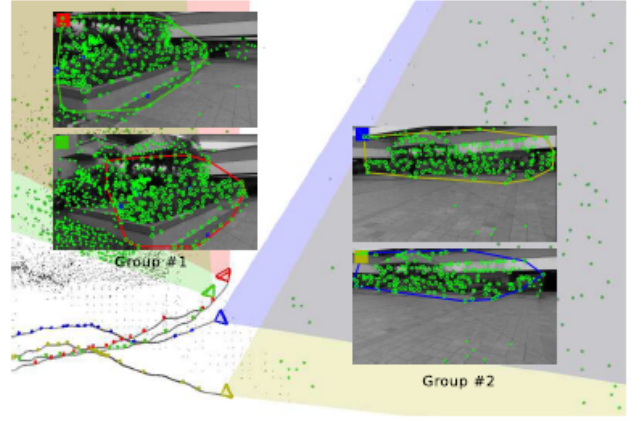


Fig. 8. The cameras were split into two groups according to their view overlap.

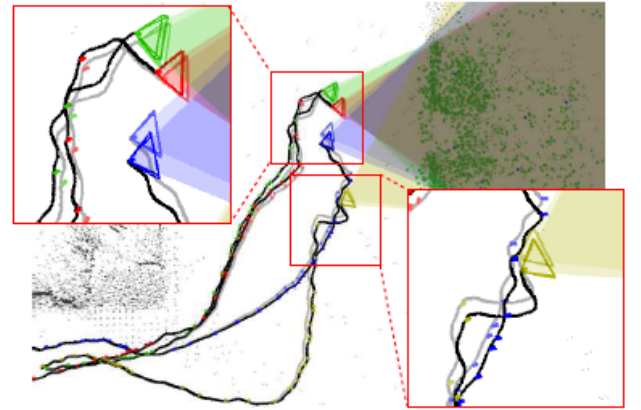


Fig. 9. The two camera groups were merged when the cameras meet again. The camera poses and map points were adjusted for consistent merge. The light gray curves (see the zoomed area) represent the camera trajectories before adjustment, and the dark ones are those after adjustment.

6.2 Merging

Two camera groups will be merged if their cameras meet and have view overlap. To detect if cameras in different groups have view overlap, we project the map points generated from one camera onto the image planes of the cameras in the other group. If the number of visible points is large ($> 30\%$ of all map points from that camera in our implementation), and the area spanned by these points are large ($> 70\%$ of the image area), we consider the two cameras to have view overlap and will merge their camera groups. A real example of such a merge is shown in Figure 9, the separated camera groups meet again.

When cameras move away from each other, the mapping and localization are performed within each camera group independently. When the cameras meet again, due to drifting errors [7], the 3D maps reconstructed from different groups are inconsistent. For example, the same object could be reconstructed at different 3D positions in different groups. Hence,

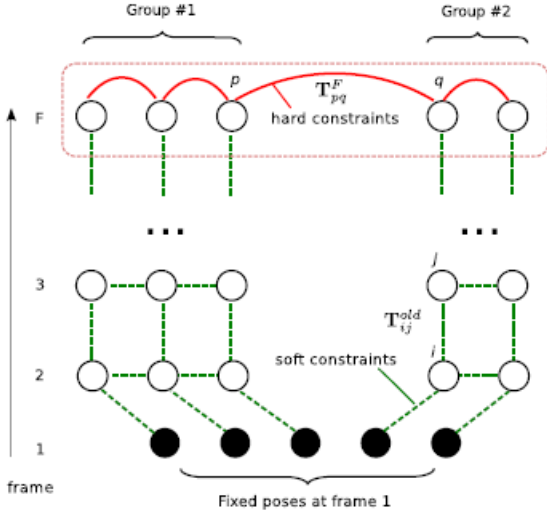


Fig. 10. Camera poses adjustment. Each vertex is a camera pose. Each edge represents a relative pose constraint, where solid and dash edges are hard and soft constraints respectively.

during group merging, we need to correct both the camera poses and map points to generate a single global consistent map. Suppose two camera groups are separated at the 1st frame and are merged at the F^{th} frame. We will adjust all camera poses from frame 2 to F , and adjust the map points generated within these frames, which consists of two successive steps described in the following section.

6.2.1 Step 1

We first estimate the correct relative poses between cameras at frame F . For this purpose, we detect and match SURF features between cameras in different groups, and then compute their relative poses (i.e. the essential matrices). We use these essential matrices to guide the matching of feature points (i.e. searching for correspondences in a narrow band within 3σ distance to the epipolar line). For each pair of matched feature points, we then merge their corresponding 3D map points by averaging their positions. In the next step, all the map points and their corresponding feature points in the F^{th} frame are put into bundle adjustment [15] to refine all camera poses.

6.2.2 Step 2

Now, we use the updated relative camera poses at the F^{th} frame as hard constraints to refine all camera poses. Figure 10 illustrates our problem formulation. We form an undirected graph where each camera pose is a vertex and each edge enforces a relative pose constraint. As shown in Figure 10, for each camera, its poses at neighboring frames are connected. For cameras in the same group, their poses at the same frame are connected if they are neighbors in the spanning tree. We fix camera poses in the 1st frame. Except the relative poses at the F^{th} frame, we treat

all the other relative poses as soft constraints. Hard and soft constraints are denoted by solid and dashed lines in Figure 10 respectively.

Let $p = 1, \dots, P$ and $q = 1, \dots, Q$ be cameras from different groups. We denote the pose of the camera p at the i^{th} frame by \mathbf{T}_p^i , and the relative pose between the camera p and q at the i^{th} frame by \mathbf{T}_{pq}^i , where

$$\mathbf{T}_p^i = \begin{pmatrix} \mathbf{R}_p^i & \mathbf{t}_p^i \\ \mathbf{0}^T & 1 \end{pmatrix}, \quad \mathbf{T}_{pq}^i = \begin{pmatrix} \mathbf{R}_{pq}^i & \alpha \mathbf{t}_{pq}^i \\ \mathbf{0}^T & 1 \end{pmatrix}. \quad (9)$$

$\mathbf{R}_p^i, \mathbf{R}_{pq}^i \in \mathbb{R}^{3 \times 3}$ and $\mathbf{t}_p^i, \mathbf{t}_{pq}^i \in \mathbb{R}^3$ are rotation matrices and translation vectors. α is used to account for the global scale difference between the two camera groups.

We treat the relative poses at the F^{th} frame as hard constraints. Hence,

$$\mathbf{T}_q^F - \mathbf{T}_{pq}^F \mathbf{T}_p^F = \mathbf{0}_{4 \times 4}, \quad (10)$$

which is equivalent to

$$\mathbf{R}_q^F - \mathbf{R}_{pq}^F \mathbf{R}_p^F = \mathbf{0}_{3 \times 3} \quad (11)$$

$$\mathbf{t}_q^F - \mathbf{R}_{pq}^F \mathbf{t}_p^F - \alpha \mathbf{t}_{pq}^F = \mathbf{0}_{3 \times 1}. \quad (12)$$

Although there are $(P + Q) \times (P + Q - 1)/2$ relative poses at the F^{th} frame, we select only $(P + Q - 1)$ of them, which either lie on the spanning trees of the camera groups or connect the two spanning trees, as illustrated by the solid lines in Figure 10. Putting all these constraints together, we get two linear systems with the following forms

$$\mathbf{U} \mathbf{r}^F = \mathbf{0} \quad \text{and} \quad \mathbf{V} \mathbf{t}^F = \mathbf{0}, \quad (13)$$

where $\mathbf{r}^F \in \mathbb{R}^{9(P+Q)}$ is a vector stacked with elements of all the rotation matrices at the F^{th} frame, and $\mathbf{t}^F \in \mathbb{R}^{3(P+Q)+1}$ is a vector that consists of all the translation elements at the F^{th} frame together with the scale factor α .

The relative camera poses from the original SLAM process are used as soft constraints. For any cameras m and n connected by the dashed edge, we expect their relative pose to have small change by the adjustment. Hence,

$$\mathbf{T}_m - \mathbf{T}_{mn}^{\text{old}} \mathbf{T}_n \approx \mathbf{0}. \quad (14)$$

Here, $\mathbf{T}_{mn}^{\text{old}}$ is the relative pose between m and n according to the SLAM process before merging. Putting all soft constraints together, we obtain two similar linear systems

$$\mathbf{A} \mathbf{r} \approx \mathbf{a} \neq \mathbf{0} \quad \text{and} \quad \mathbf{B} \mathbf{t} \approx \mathbf{b} \neq \mathbf{0}, \quad (15)$$

where $\mathbf{r} \in \mathbb{R}^{9F(P+Q)}$ and $\mathbf{t} \in \mathbb{R}^{3F(P+Q)}$ are vectors stacked by all the rotation and translation elements of all frames. Notice that the right sides of the two linear systems are not equal to zero because the camera poses at the 1st frame are fixed.

Combining both the hard constraints in Equation (10) and soft constraints in Equation (14), we obtain

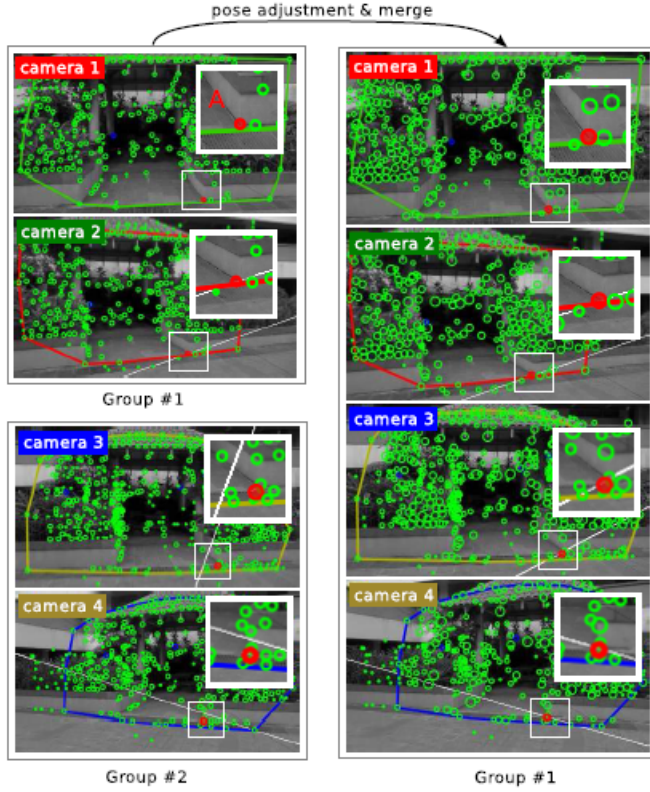


Fig. 11. Camera poses and map point positions are adjusted during group merge. Left: before the adjustment, the corresponding feature points of A in the ‘camera 3’ and ‘camera 4’ do not lie on the its epipolar lines. Right: after adjustment, all its corresponding feature points lie on the epipolar lines.

the updated cameras poses and the scale factor by solving two constrained linear least square problems

$$\arg \min_r \|Ar - a\|^2 \quad \text{s.t.} \quad \hat{U}r = 0 \quad (16)$$

and

$$\arg \min_{\hat{t}} \|\hat{B}\hat{t} - \hat{b}\|^2 \quad \text{s.t.} \quad \hat{V}\hat{t} = 0, \quad (17)$$

where $\hat{t} \in \mathbb{R}^{3F(P+Q)+1}$ is t appended with a scale factor α . $\hat{U}, \hat{V}, \hat{B}, \hat{b}$ are the augmented matrices and vectors by adding zero elements. Note that we do not impose orthonormality condition to the rotation matrices in this formulation. Hence, once we obtain results from the above two equations, we further find the closest rotation matrices to the initial matrices by SVD (i.e. setting all the singular values to one).

The above optimization problem is converted to a set of sparse linear equations [13]. We use the CSparse[8] library to solve them in our system. After the camera poses have been updated, the 3D positions of map points are also updated by re-triangulating their corresponding feature points.

In Figure 9, the camera poses are updated when the two camera groups merge. The light gray curves (and dark curves) represent the camera trajectories before (and after) merging. To further exemplify the

importance of pose updates, we examine the epipolar geometries in Figure 11. In the left of Figure 11, we plot the epipolar lines of the feature point A in the first camera. Because of the map inconsistency between the two camera groups, the corresponding feature points in the third and the fourth camera do not lie on the epipolar lines. In comparison, after camera pose update and re-triangulation of map points, the corresponding feature points in all cameras lie on their respective epipolar lines as shown in the right of Figure 11. To provide an additional validation, when visualizing the projection of a map point in Figure 11, we use its circle radius to indicate its number corresponding feature points. It is clear that the overall circle size is larger on the right, which suggests map points have more corresponding feature points after the adjustment. This increase in corresponding feature points comes from the merge of duplicated 3D map points.

7 INCREMENTAL REFINEMENT

We refine both the camera poses and the 3D map points from time to time by bundle adjustment. For better efficiency, bundle adjustment only refines the camera poses of some selected key frames and the map points reconstructed from these frames. Whenever there is a significant drop (30%) in the number of tracked feature points in any camera, we insert a key frame for all cameras.

The bundle adjustment runs in a separate thread, which operates with the most recent K key frames. It is called when $K - 1$ key frames have been inserted consecutively (i.e. in two successive bundle adjustment calls, there is one common key frame). The bundle adjustment only refines camera poses of key frames and map points reconstructed from these frames. To refine the camera poses of the other frames, we adopt a similar method as Section 6.2.2. Basically, we fix the camera pose at key frames, and use the relative poses between successive frames before bundle adjustment as soft constraint. In other words, we enforce $T_m - T_{mn}^{old} T_n \approx 0$, where T_{mn}^{old} is the relative pose between camera m, n before bundle adjustment. We then update all the camera poses while keep those at key frames unchanged. After the pose refinement, the 3D positions of other map points are updated by re-triangulating their corresponding feature points.

8 RESULTS

We tested our collaborative SLAM system on both static scenes and dynamic scenes. All data were captured by hand held cameras with Field of View about 70° and processed offline. In all our experiments, we set the standard deviation of feature detection uncertainty σ as 3.0 pixels. The threshold for Mahalanobis distance θ is set to be 2.0 to decide if a feature point is an inlier or outlier (with confidence of 95% according

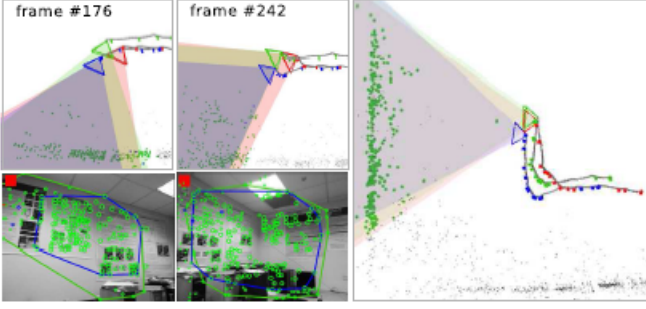


Fig. 12. Our results on the ‘wall’ sequence. Trajectories of different cameras are visualized in different color from the top-down viewpoint. Map points are projected in the input video frames for a reference.

to Gaussian distribution). The ZNCC threshold τ_{ncc} to measure the similarity between image patches is set to 0.7. The minimum number of frames N_{min} to triangulate a feature track is set to 60. The number of frames N_{rec} for active map point caching is 200. The radius ϕ_r for searching nearby seed matches in ‘inter-camera mapping’ is set to 10% of $\max\{\text{image width, image height}\}$, and $\phi_d = 3\phi_r$. In practice, we found our results are not sensitive to these parameters. We also present the results on an accompanying video.

8.1 Static Scenes

8.1.1 Critical camera motion

Single camera based visual SLAM systems usually fail under critical motion, such as camera rotation without translation. This is because when a camera rotates, the number of visible map points drops quickly. However, new map points cannot be generated because of little camera translation. This problem is more serious when the camera field of view is small. Our CoSLAM system can deal with such situation by the collaboration among multiple cameras as demonstrated in ‘wall’ example in Figure 12. In this example, cameras started rotating at around the 180th frame and finished at about the 256th frame. Our method successfully captured the motion of all cameras and the two perpendicular walls.

For a comparison, we applied the single camera based SLAM system, PTAM [18], on the same data. The results from PTAM are provided in Figure 13. PTAM failed on all cameras when they started to rotate. Even the re-localization cannot recover it from this failure. This is because the re-localization works only when the newly captured frames have view overlap with previous key frames. In camera rotation, the camera is turning to a novel view that is not observed before. Hence, re-localization cannot help in this case.

8.1.2 Drift analysis

We tested our CoSLAM system in a middle scale ‘courtyard’ example. The average length of camera

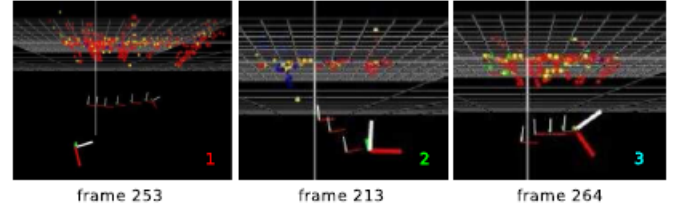


Fig. 13. Results generated by PTAM [18] on the ‘wall’ sequence. The system failed on all three sequences when the camera started rotating.

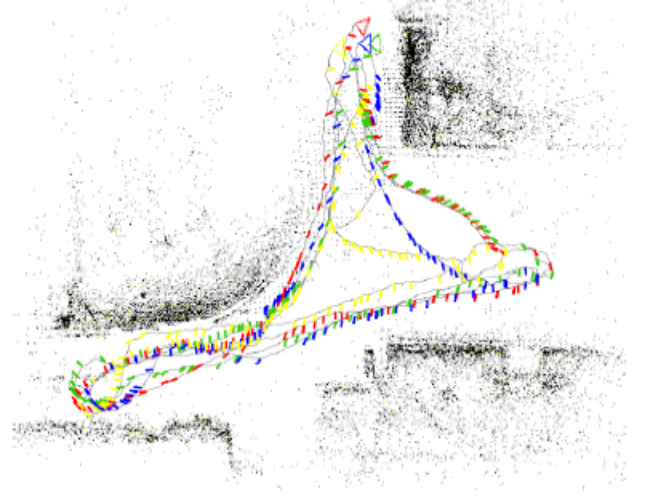


Fig. 14. The overview of the ‘courtyard’ example. Four hand held cameras were used to capture the data. The view overlaps between cameras changed over time.

trajectories was 96 m. The input videos were captured in a courtyard by four hand held cameras. The view overlap among cameras changed over time, which led to camera group splitting and merging as discussed in Section 6. During data capturing, we walked around the courtyard and returned to the starting place. Hence, the drift error can be measured by manually identifying corresponding points in the first and the last video frames. We analyzed the drift errors with different number of cameras in the system. An overview of the scene is provided in Figure 14.

We tried all possible 1-camera, 2-camera, 3-camera, and 4-camera CoSLAM. The average distance drift errors were 2.53m, 1.57m, 1.19m, and 0.67m respectively. The average scale drift error were 0.76, 1.10, 0.96, 1.00 respectively. This result is visualized in Figure 15, where the red and green line segment indicates the reconstruction of the sign board from the first and last frames. This result indicates that CoSLAM with multiple cameras can successfully reduce the drift errors.

For a comparison, we tried to apply the PTAM system to this example and measure its drift error. However, PTAM failed at all the four 1-camera tests and cannot finish the whole loop. In comparison, our system only failed in one 1-camera test. We believe this is because we maintain the triangulation uncer-

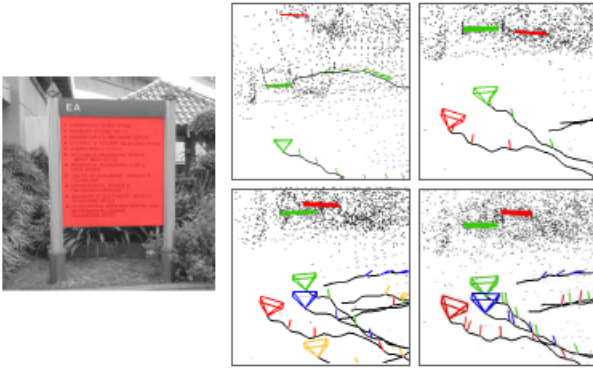


Fig. 15. The final drift error is measured by the difference between the board reconstructed from the first (marked by red) and the last video frames (marked by green). The left is an example image of the board. The right are the CoSLAM results with different number of cameras.

tainty of map points, and continuously improve map accuracy when more images are captured.

8.2 Dynamic Scenes

The capability of our method to robustly estimate camera poses in highly dynamic scenes is demonstrated in Figure 16. Note that the number of static points (green dots) was small in all cameras. Further, the static points were often distributed within a small region, which usually leads to large error if the camera pose is estimated only according to static points. Our CoSLAM system automatically switched to the ‘inter-camera pose estimation’, and successfully estimated the poses of all cameras in such a challenging case. For a comparison, we disabled the ‘inter-camera pose estimation’ and applied the ‘intra-camera pose estimation’ to the same sequence. The system failed at the # 665-th frame because of the large error in pose estimation. This comparison can be seen clearly from the visualization of the camera trajectories in the right column of Figure 16.

We tested our CoSLAM system in an indoor scene, the ‘walking man’ example in Figure 17. The relatively dim indoor lighting usually leads to blurry images, which make feature tracking difficult and pose estimation inaccurate. Our CoSLAM system can successfully handle this data. The estimated map points and camera trajectories are visualized from the top view in the left of Figure 17. An important feature of our system is that we can recover the 3D trajectories of moving points, which is demonstrated in Figure 18, where the blue curves are 3D trajectories of the dynamic points on the walking person. We also manually specify corresponding points to measure the drift error. Our CoSLAM system had 1.2m distance drift and 1.12 scale drift. The average length of camera trajectories is 28.7m.

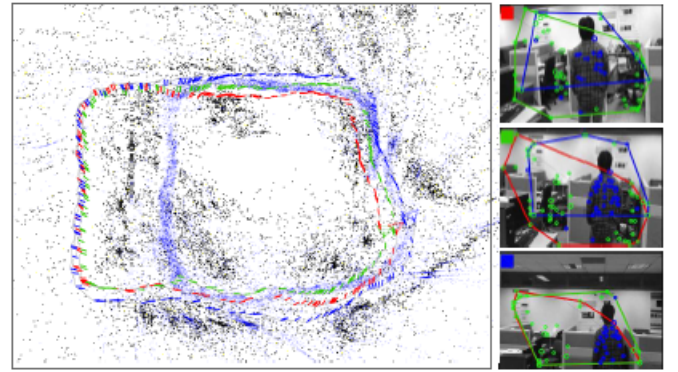


Fig. 17. Result on the ‘walking man’ example. The blue points represent the moving points in the scene.

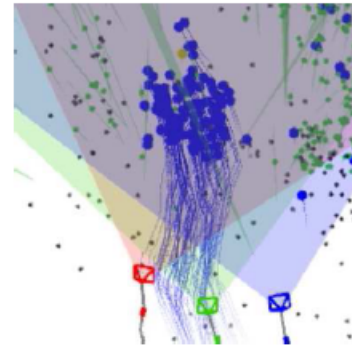


Fig. 18. Our system can track the time-varying 3D positions of the dynamic points on the moving objects.

We further tested the robustness of our system on a middle scale ‘garden’ example, where the average length of camera trajectories is 63m. Three hand held cameras were used to capture the input videos in a garden. Several people walked in front of the camera to disturb the capturing process. As the moving people frequently occupied a large part of the video frame, the SLAM problem with this data was very challenging (please refer to the supplementary video). As shown in Figure 19, our method succeeded in such a data. The manually measured drift errors are 5.3m in distance and 0.65 in scale. These errors are relative large, because the moving objects frequently occluded nearby static points. In this situation, the faraway static points played a more important role in camera pose estimation. However, the positions of these faraway points were less reliable, which led to inaccurate estimation and finally produced a relative large drift error.

8.3 Run Time Efficiency

Our system was developed under the Ubuntu 64-bit system with an Intel i7 CPU (4 cores at 2.80GHz), 4G RAM and an nVidia GeForce GTX 480 graphics card. The main thread of the system estimated the egomotion of the cameras and generated map points. Bundle adjustment and camera group management

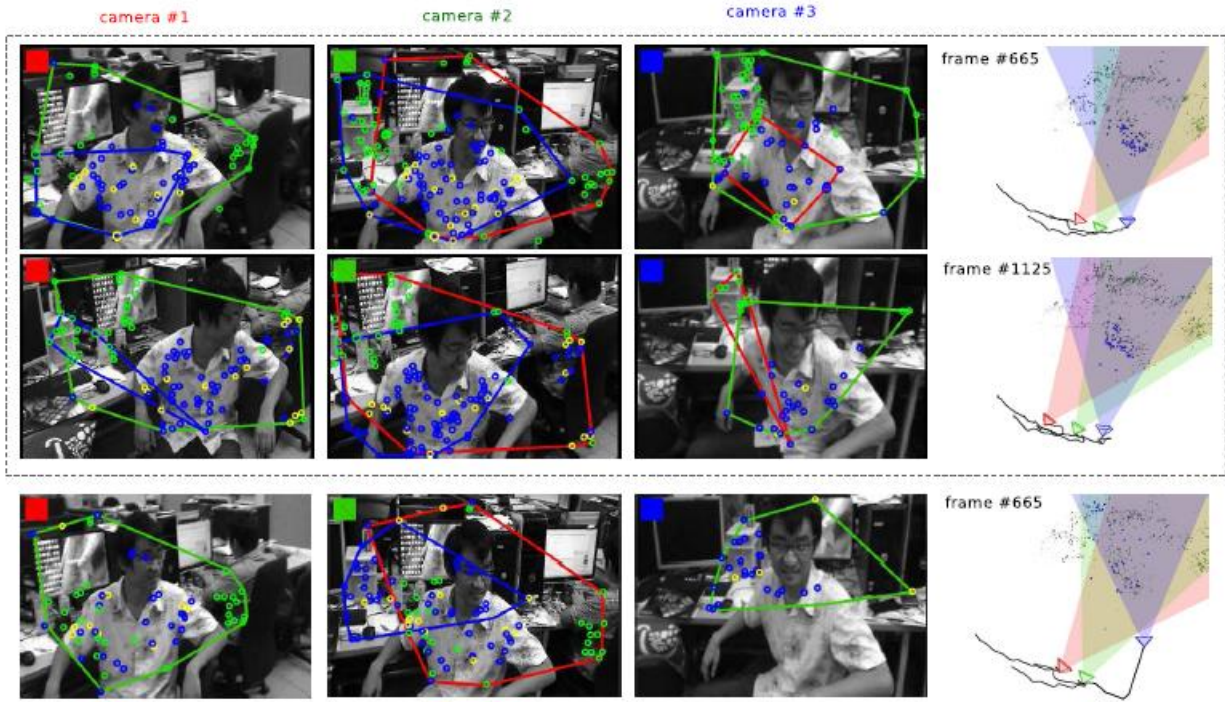


Fig. 16. Our results on the ‘sitting man’ sequence. There are only a few static points in the scene. Our ‘inter-camera pose estimation’ successfully estimates all camera poses (shown in the first two rows). The bottom row shows the result by only applying the ‘intra-camera pose estimation’. The pose of the blue camera (camera #3) was completely wrong at the 665-th frame, because of little static points. The reconstructed map points (both static and dynamic points) and camera trajectories of the last 150 frames are provided on the right (visualized from the top view)

TABLE 1
Average timings

components	ms	calling conditions
feature tracking	8.7	every frame (by GPU)
intra-camera pose estimation	10.7	every frame
inter-camera pose estimation	57.2	see Section 4
map point classification	4.9	every frame
map point registration	14.47	every frame
intra-camera mapping	2.3	see Section 5.2
inter-camera mapping	48.3	see Section 5.2

were called in separated threads. Feature tracking was implemented in GPU according to the GPU-KLT [29].

We evaluated the run time efficiency on the ‘garden’ sequence shown in Figure 19. All the 4800 frames were used for evaluation. The average time spent in each call of all components are listed in Table 1. Most of the components run quickly. The ‘inter-camera pose estimation’ and the ‘inter-camera mapping’ took about 50 ms to run. The number of map points (including both dynamic and static points) and the processing time of each frame are also shown respectively in Figure 21. Although the runtime efficiency was reduced when inter-camera operations were called (see the peaks in Figure 21), our system on average ran in real-time and took about 38ms to process a frame with about one thousand of map points.

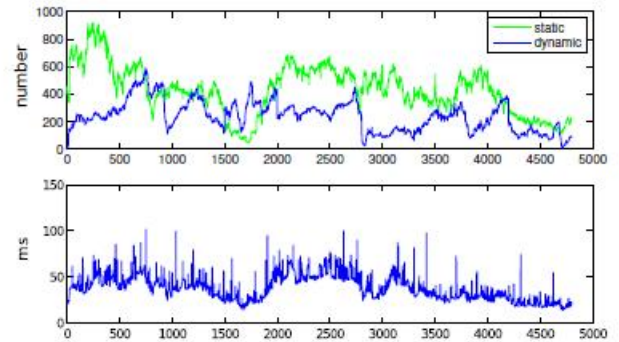


Fig. 21. Run time efficiency with three cameras. Top: number of map points over time. Bottom: the time spent to process each frame. The average time to process a frame is 38ms.

8.4 System Scalability

We tested our system scalability with 12 cameras moving independently in a static scene. Some results are shown in Figure 20. On the left of each row shows the top view of the reconstructed scene map with camera trajectories. On the right are the input images from all 12 cameras. Images framed by the same color come from cameras in the same group. This scene contains several separated ‘branches’ as can be seen from the scene map in the third row. These 12 cameras were divided into several troops, and each

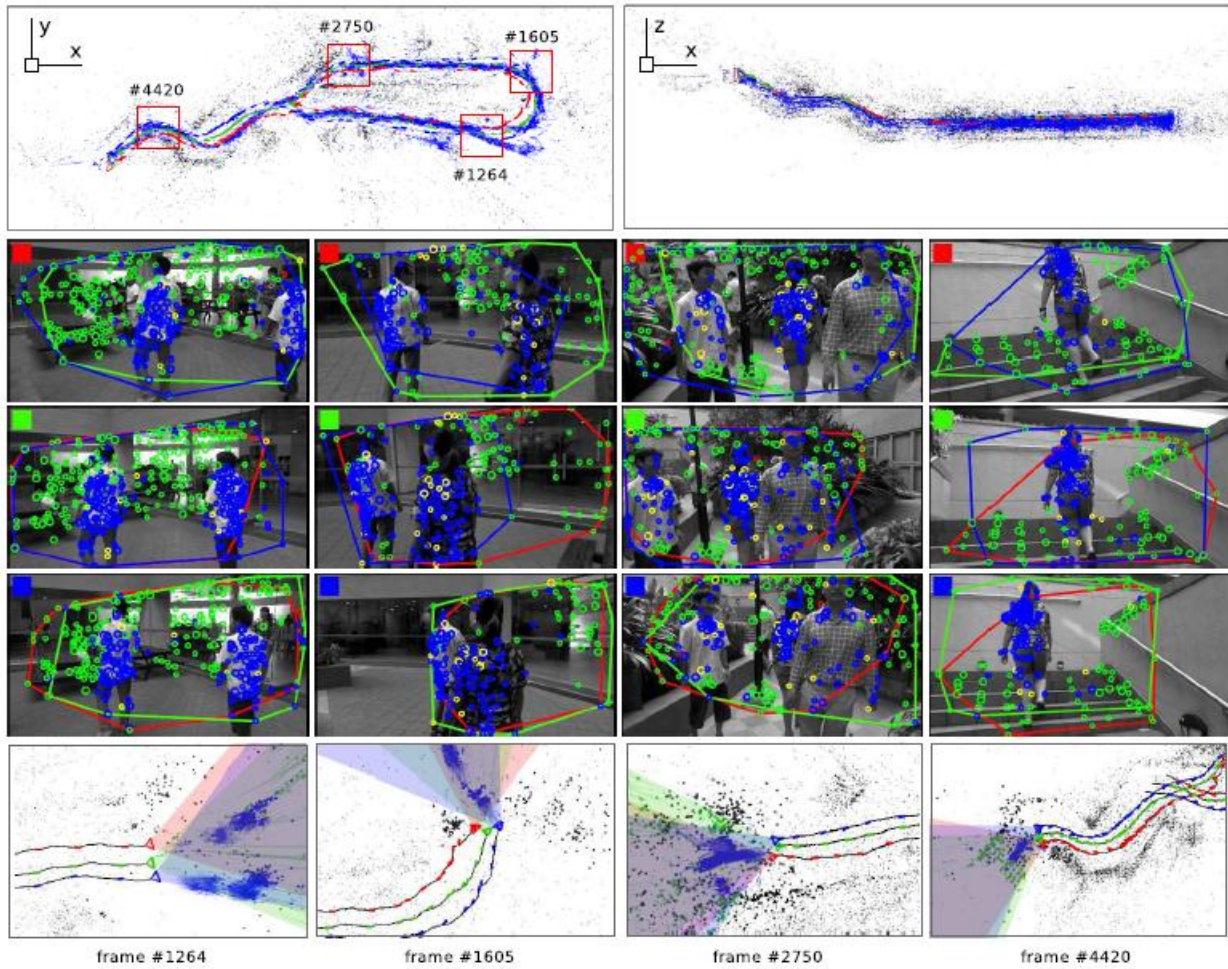


Fig. 19. The collaborative SLAM result in a challenge dynamic scene using three cameras. An overview of the reconstruction is provided in the first row from the top and side view respectively. The gray point cloud indicates the static scene structures, while the blue trajectories are the moving points. In the middle are the detected feature points on video frames. The green and blue points represent static and dynamic points respectively, while the yellow points are the unmapped feature points (please zoom in to check details). In the last row we provide some zoomed views of the camera trajectories with their frame index marked in the first row.

troop explored one branch. Our system can correctly handle camera group splitting and merging. (Please refer to the supplementary video.) However, when there are 12 cameras, the average run time efficiency dropped significantly to about 1fps due to the heavy computation.

9 CONCLUSION AND FUTURE WORK

We propose a novel collaborative SLAM system with multiple moving cameras in a possibly dynamic environment. The cameras move independently and can be mounted on different platforms, which makes our system potentially applicable to robot teams [4], [34], and wearable augmented reality [5]. We address several issues in pose estimation, mapping and camera group management, so that the system can work robustly in challenging dynamic scenes as shown in the experiments. The whole system runs in real-time. Currently, our system works offline with pre-captured

video data. We plan to integrate a data capturing component to make an online system. Further, our current system requires synchronized cameras, and all images from these cameras are sent back and processed in the same computer. It will be interesting to develop a distributed system for collaborated SLAM, where computation is distributed to multiple computers.

10 ACKNOWLEDGEMENT

This work is supported by the Singapore grant R-263-000-555-112, R-263-000-620-112 and AORAD grant R-263-000-673-597.

REFERENCES

- [1] J. Allred, A. Hasan, S. Panichsakul, W. Pisano, P. Gray, J. Huang, R. Han, D. Lawrence, and K. Mohseni. Sensorflock: an airborne wireless sensor network of micro-air vehicles. In *Proc. of Int'l Conf. on Embedded networked sensor systems*, pages 117–129. ACM, 2007.
- [2] C. Bibby and I. Reid. Simultaneous localisation and mapping in dynamic environments (slamde) with reversible data association. In *Proc. of Robotics: Science and Systems*, 2007.

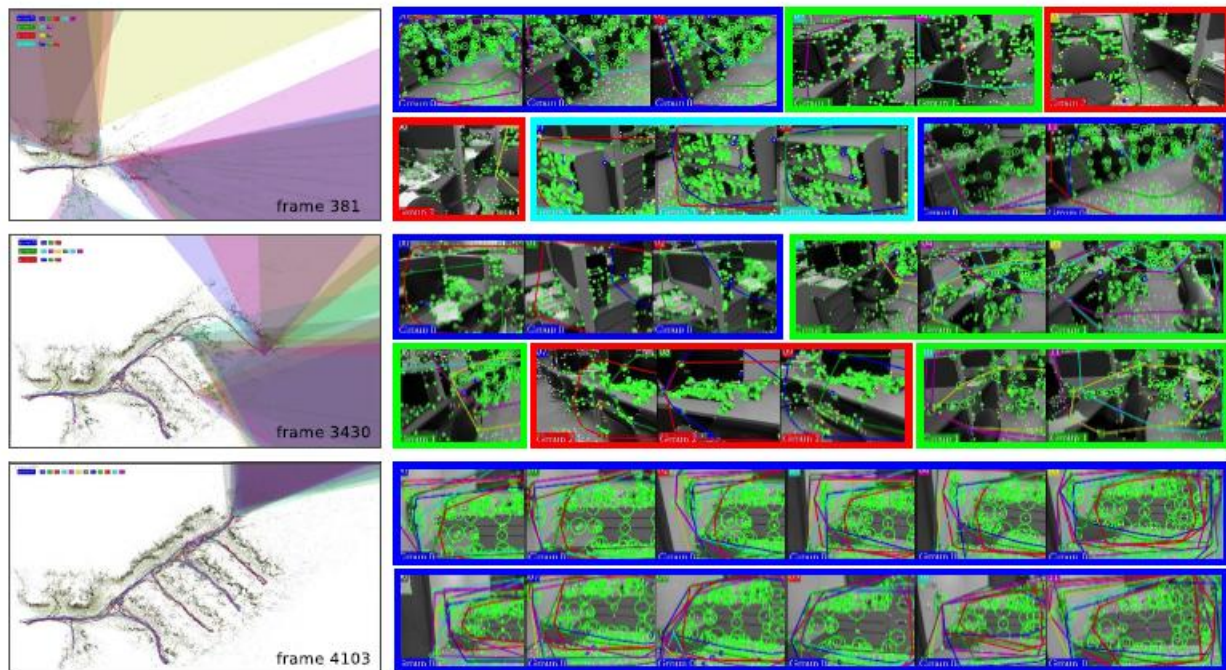
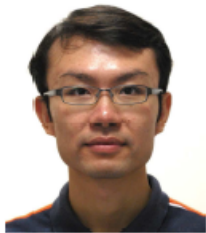


Fig. 20. Our system was tested with 12 cameras. The map and camera poses are shown in the left, where the top left corner shows the camera grouping. Images and feature points are shown in the right, where images framed in the same color come from cameras in the same group. These three rows have four, three and one camera groups respectively.

- [3] C. Bibby and I. Reid. A hybrid slam representation for dynamic marine environments. In *ICRA*, pages 257–264. IEEE, 2010.
- [4] W. Burgard, M. Moors, D. Fox, R. Simmons, and S. Thrun. Collaborative multi-robot exploration. In *IEEE Proc. of Robotics and Automation*, volume 1, pages 476–481, 2002.
- [5] R. O. Castle, G. Klein, and D. W. Murray. Video-rate localization in multiple maps for wearable augmented reality. In *Proc 12th IEEE Int Symp on Wearable Computers*, pages 15–22, 2008.
- [6] M. Chli and A. J. Davison. Active matching for visual tracking. *Robot. Auton. Syst.*, 57(12):1173–1187, 2009.
- [7] K. Cornelis, F. Verbiest, and L. Van Gool. Drift detection and removal for sequential structure from motion algorithms. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 26(10):1249–1259, 2004.
- [8] T. Davis. *Direct Methods for Sparse Linear Systems*. SIAM, 2006.
- [9] A. Davison. Real-time simultaneous localisation and mapping with a single camera. In *IEEE Proc. of ICCV*, pages 1403–1410, 2003.
- [10] A. Davison, I. Reid, N. Molton, and O. Stasse. MonoSLAM: Real-time single camera SLAM. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, pages 1052–1067, 2007.
- [11] H. Durrant-Whyte and T. Bailey. Simultaneous localisation and mapping (slam): Part i the essential algorithms. In *IEEE Robotics and Automation Magazine*. Citeseer, 2006.
- [12] E. Eade and T. Drummond. Scalable monocular SLAM. In *IEEE Proc. of CVPR*, volume 1, pages 469–476, 2006.
- [13] G. Golub. Numerical methods for solving linear least squares problems. *Numerische Mathematik*, 7(3):206–216, 1965.
- [14] D. Hahnel, R. Triebel, W. Burgard, and S. Thrun. Map building with mobile robots in dynamic environments. In *IEEE Proc. of Robotics and Automation*, volume 2, pages 1557–1563, 2003.
- [15] R. Hartley and A. Zisserman. *Multiple view geometry*, volume 6. Cambridge university press, 2000.
- [16] K. Ho and P. Newman. Detecting loop closure with scene sequences. *Int'l Journal of Computer Vision*, 74(3):261–286, 2007.
- [17] M. Kaess and F. Dellaert. Visual slam with a multi-camera rig. *Georgia Institute of Technology, Tech. Rep. GIT-GVU-06-06*, 2006.
- [18] G. Klein and D. Murray. Parallel tracking and mapping for small AR workspaces. In *IEEE & ACM Proc. of Int'l Symp. on Mixed and Augmented Reality*, pages 225–234, 2007.
- [19] A. Kundu, K. Krishna, and C. Jawahar. Realtime motion segmentation based multibody visual slam. pages 251–258, 2010.
- [20] B. Leibe, N. Cornelis, K. Cornelis, and L. Van-Gool. Dynamic 3d scene analysis from a moving vehicle. In *IEEE Proc. of CVPR*, 2007.
- [21] E. Mouragnon, M. Lhuillier, M. Dhome, F. Dekeyser, and P. Sayd. Real time localization and 3d reconstruction. In *IEEE Proc. of CVPR*, volume 1, pages 363–370, 2006.
- [22] R. Newcombe and A. Davison. Live dense reconstruction with a single moving camera. In *IEEE Proc. of CVPR*, pages 1498–1505, 2010.
- [23] D. Nister, O. Naroditsky, and J. Bergen. Visual odometry. In *IEEE Proc. of CVPR*, volume 1, 2004.
- [24] K. Ozden, K. Schindler, and L. V. Gool. Multibody structure-from-motion in practice. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, pages 1134–1141, 2010.
- [25] L. Paz, P. Piniés, J. Tardós, and J. Neira. Large-scale 6-dof slam with stereo-in-hand. *IEEE Trans. on Robotics*, 24(5):946–957, 2008.
- [26] E. Royer, M. Lhuillier, M. Dhome, and T. Chateau. Localization in urban environments: monocular vision compared to a differential gps sensor. 2:114–121, 2005.
- [27] E. Sahin. Swarm robotics: From sources of inspiration to domains of application. *Swarm Robotics*, pages 10–20, 2005.
- [28] J. Shi and C. Tomasi. Good features to track. In *IEEE Proc. of CVPR*, pages 593–600, 1994.
- [29] S. Sinha. http://www.cs.unc.edu/~ssinha/Research/GPU_KLT/.
- [30] P. Smith, I. Reid, and A. Davison. Real-time monocular SLAM with straight lines. In *Proc. British Machine Vision Conference*, volume 1, pages 17–26, 2006.
- [31] H. Strasdat, J. Montiel, and A. Davison. Real-time monocular SLAM: Why filter? In *IEEE Proc. of Robotics and Automation*, pages 2657–2664, 2010.
- [32] H. Strasdat, J. Montiel, and A. Davison. Scale drift-aware large scale monocular slam. In *Proc. of Robotics: Science and Systems*, 2010.
- [33] H. Strasdat, J. Montiel, and A. Davison. Visual SLAM: Why Filter? *Image and Vision Computing*, 2012.
- [34] S. Thrun, W. Burgard, and D. Fox. A real-time algorithm for mobile robot mapping with applications to multi-robot and 3D mapping. In *IEEE Proc. of Robotics and Automation*, volume 1, pages 321–328. IEEE, 2002.
- [35] C. Wang, C. Thorpe, S. Thrun, M. Hebert, and H. Durrant-

- Whyte. Simultaneous localization, mapping and moving object tracking. *Int'l of Robotics Research*, 26(9):889, 2007.
- [36] B. Williams, G. Klein, and I. Reid. Real-Time SLAM Relocalization. In *IEEE Proc. of ICCV*, pages 1–8, 2007.
- [37] N. Winters, J. Gaspar, G. Lacey, and J. Santos-Victor. Omnidirectional vision for robot navigation. In *IEEE Work. on Omnidirectional Vision*, page 21, 2000.
- [38] D. Wolf and G. Sukhatme. Mobile robot simultaneous localization and mapping in dynamic environments. *Autonomous Robots*, 19(1):53–65, 2005.
- [39] Z. Zhang. Parameter estimation techniques: A tutorial with application to conic fitting. *Image and vision Computing*, 15(1):59–76, 1997.
- [40] J. ZUFFEREY. *Bio-Inspired vision-based flying robots*. PhD thesis, ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE, 2005.



Danping Zou received the BS degree from Huazhong University of Science and Technology (HUST) in 2003 and the PhD degree from Fudan University in 2010 in China. He is now a research fellow in the Department of Electrical and Computer Engineering at National University of Singapore. His research interests include video tracking and low-level 3D vision.



Ping Tan received the BS degree in Applied Mathematics from the Shanghai Jiao Tong University in China in 2000 and the PhD degree in Computer Science and Engineering from the Hong Kong University of Science and Technology in 2007. He joined the Department of Electrical and Computer Engineering at the National University of Singapore as an assistant professor in 2007. He received the MIT TR35@Singapore award in 2012. His research interests include computer vision and computer graphics. He is a member of the IEEE and ACM.