

Collaborative Dense SLAM

Louis Gallagher and John B. McDonald

Department of Computer Science, Maynooth University, Maynooth, Co. Kildare, Ireland

louis.gallagher.2013@mumail.ie, johnmcd@cs.nuim.ie

Abstract

In this paper, we present a new system for live collaborative dense surface reconstruction. Cooperative robotics, multi participant augmented reality and human-robot interaction are all examples of situations where collaborative mapping can be leveraged for greater agent autonomy. Our system builds on ElasticFusion to allow a number of cameras starting with unknown initial relative positions to maintain local maps utilising the original algorithm. Carrying out visual place recognition across these local maps the system can identify when two maps overlap in space, providing an inter-map constraint from which the system can derive the relative poses of the two maps. Using these resulting pose constraints, our system performs map merging, allowing multiple cameras to fuse their measurements into a single *shared* reconstruction. The advantage of this approach is that it avoids replication of structures subsequent to loop closures, where multiple cameras traverse the same regions of the environment. Furthermore, it allows cameras to directly exploit and update regions of the environment previously mapped by other cameras within the system.

We provide both quantitative and qualitative analyses using the synthetic ICL-NUIM dataset and the real-world Freiburg dataset including the impact of multi-camera mapping on surface reconstruction accuracy, camera pose estimation accuracy and overall processing time. We also include qualitative results in the form of sample reconstructions of room sized environments with up to 3 cameras undergoing intersecting and loopy trajectories.

Keywords: 3D reconstruction, Dense Mapping, Collaborative Mapping, SLAM, Machine Vision

1 Introduction

Simultaneous localisation and mapping (SLAM) is a core problem in robotics and computer vision, providing the foundations for many higher level environment understanding and interaction tasks, particularly in autonomous mobile robotics and augmented reality. With the advent of consumer grade active depth sensors, such as the Microsoft Kinect and the Asus Xtion Pro, and the additional compute power made available through modern GPUs, dense mapping has become a prominent focus of research in visual SLAM. Recently, the state-of-the-art in the field has come to include a multitude of systems offering large-scale, high-precision mapping and tracking capabilities [Dai et al., 2017, Whelan et al., 2016, Kerl et al., 2013, Whelan et al., 2014, Newcombe et al., 2011a, Engel et al., 2014, Mur-Artal and Tardos, 2017, Newcombe et al., 2011b, Keller et al., 2013]. What distinguishes dense mapping from its sparse counterpart is twofold; firstly, a direct approach to camera tracking where potentially every pixel is used; secondly, the estimation of a map that consists of a dense set of measurements of the surface that underlies the scene. Examples of such map representations include volumetric signed distance functions and surface element (surfel) lists.

Much of the developments within dense SLAM have been focused on scenarios involving mapping the environment using a single-sensor platform. However, many scenarios exist where multiple mobile agents

This research was supported, in part, by the IRC GOIPG scholarship scheme grant GOIPG/2016/1320 and, in part, by Science Foundation Ireland grant 13/RC/2094 to Lero - the Irish Software Research Centre (www.lero.ie)

must work together in the same space, thereby requiring a single shared map of the environment. The main contribution of this paper is a system that addresses the wider problem of online collaborative, multi-camera dense mapping and tracking using RGB-D sensors with unknown initial relative poses.

Historically collaborative mapping has been a recurring theme in the SLAM literature where many problems that emerge when multiple cameras must map collaboratively have been addressed [Saeedi et al., 2016, Fenwick et al., 2002, Williams et al., 2002, Thrun and Liu, 2003, Howard, 2006, Howard et al., 2006, Zhou and Roumeliotis, 2006, McDonald et al., 2013]. For example, in general the starting position of cameras relative to one another is not known and so they do not share a common reference frame, which instead must be identified online, as mapping proceeds. Once a common reference frame has been found maps must be merged, from which point updates to the map from multiple cameras must be coordinated so as to avoid corrupting the map.

Collaborative mapping has come back into focus in SLAM research in recent times. Golodetz et al. propose a novel approach to dense collaborative mapping that uses InfiniTAM [Prisacariu et al., 2017] to construct *separate* voxel-based maps for each camera. The system uses a random forest relocaliser to find estimates of the relative transforms between each of these camera specific maps. These relative transforms act as constraints in a pose graph which, once optimised, yields a set of optimal global poses, one for each of the sub maps. Although these global poses align the submaps into a common reference frame they are never merged directly, instead each camera maintains its own submap throughout [Golodetz et al., 2018]. In contrast, our system performs map merging, allowing multiple cameras to fuse their measurements into a single *shared* reconstruction. The advantage of this approach is that it avoids replication of structures subsequent to loop closures, where multiple cameras traverse the same regions of the environment. Furthermore, it allows cameras to directly exploit and update regions of the environment previously mapped by other cameras within the system.

The remainder of this paper is structured as follows. In Section 2 we briefly review the ElasticFusion (EF) dense mapping system. In Section 3 we describe how we extended EF to allow multi-camera mapping. In section 3.1 we explain the inter-map loop closure mechanism and how maps are merged. Later, in Section 4 we describe the experimental process used to measure the systems performance and present the results of this experimentation. Finally, in Section 5 we conclude with some remarks on the proposed system and give some directions for future research.

2 Background

In this section we briefly summarise the main elements of the ElasticFusion (EF) algorithm and define the notation used for the remainder of the paper. For a more comprehensive treatment of EF see [Whelan et al., 2016]. In EF the map is represented as an unordered list of surface elements (surfels), \mathcal{M} . Each surfel is an estimate of a discrete patch centred around an associated scene point. A surfel is characterised by its position $p \in \mathbb{R}^3$, normal $n \in \mathbb{R}^3$, radius $r \in \mathbb{R}$, weighting $w \in \mathbb{R}$, colour $c \in \mathbb{N}^3$, initialisation time t_0 and the time it was last observed t . The weight of a surfel signifies the level of confidence in the estimation of its parameters. \mathcal{M} is split along temporal lines into two non-intersecting sub-lists, Θ containing surfels that are considered active, meaning they have recently been observed by the camera, and Φ containing those surfels that have not been observed recently and are considered inactive.

Camera tracking is achieved by aligning each input frame at time t with a predicted view of the model which is rendered from Θ using the pose P_{t-1} . The transformation that brings the two frames into alignment is estimated by minimising a combined photometric and geometric error over a three level Gaussian pyramid. The resulting transformation is applied to P_{t-1} to yield an estimate of the global pose of the camera, P_t , for the current time step.

Following camera tracking, in the next stage of the pipeline EF seeks to close two types of loop; global and local loops. Global loop closures are identified with a fern-based visual place recognition system as proposed by [Glocker et al., 2015]. Using the pose P_t , from the camera tracking step, a new active model prediction is rendered from which a fern encoding is computed and used to search the fern database for a matching view. If one is found then the two views are aligned, generating a set of surface-to-surface constraints which are used to optimise the nodes in a deformation graph. The nodes of this graph are initialised by subsampling a set of

surfels from \mathcal{M} . Once optimised, the graph can be applied directly to the surface. To do this an influence region is computed for each surfel consisting of the set of graph nodes close both spatially and in initialisation time. The final pose of a surfel is given by the weighted application of the affine transformations held in each of the nodes in the surfel’s influence region where the weight of a given node is a function of its distance from the surfel.

If a global loop closure has not occurred at the current time step then local loop closures aligning Θ to Φ are sought. Local loops are closed by performing a surface registration between the portions of Θ and Φ in view of P_t . Similar to global loop closures this registration yields a set of surface-to-surface constraints that are used to optimise the nodes of a deformation graph. Once the deformation graph has been applied the surfels in Φ that were in view are reactivated within Θ .

Fusion of the current frame can proceed once the camera has been tracked and all loop closures have been resolved. Surfels in Θ are projectively associated with pixels in the current frame using P_t and K , the camera’s intrinsic parameters. The normals, vertices, colours and radius of associated points are merged with a weighted average of the model surfel and the newly measured vertex. The weight of the surfel is updated in accordance with the normalised radial distance of the measured vertex from the camera centre. If there is no corresponding surfel for a vertex then a new unstable surfel is inserted into the model. Surfels become stable after their confidence counter reaches a threshold through repeated measurement, from which point they can be used for camera tracking.

3 Approach

In this section we extend the EF system to permit dense mapping using multiple independently moving RGB-D cameras to map a space in a collaborative fashion. As the cameras explore a space they fuse their measurements into, and track against, a shared model of the reconstruction.

In extending EF to permit collaborative mapping we identify two distinct phases of processing for any given input stream. In the initial phase, the system assumes that each camera is positioned at the origin of a frame of reference that is independent to that of other cameras essentially processing it via an independent EF mapping pipeline. As mapping progresses the aim is for inter-camera loop closures to occur, thereby providing the necessary transformations between the camera submaps. These transformations permit alignment of submaps into a common frame of reference which makes it possible to fuse subsequent measurements from each camera into a single global map.

A set of n cameras $\{\mathbb{C}_i\}, i \in [1..n]$ serve frames to a central server that runs our mapping process. Each \mathbb{C}_i is represented by a 4×4 transformation matrix from the Lie group \mathbb{SE}_3 . We use the Lightweight Communications and Marshalling (LCM) system of [Huang et al., 2010] to connect the cameras to the central mapping server over a network. LCM provides efficient encoding and decoding mechanisms for transmission of camera data over a network and the subsequent unpacking of that data into data structures that are more amenable to processing.

Given that there are multiple cameras each frame is of the form \mathbb{F}_i^t where the superscript t indicates the frame’s timestamp and the subscript i denotes that this frame is from camera \mathbb{C}_i . We note that the cameras need not be synchronised in time.

If we assume momentarily that all cameras are globally aligned then multi-camera fusion proceeds in essentially the same way as in the original EF algorithm. At each timestep t the frames from each camera for that timestep are processed in *sequence* by EF. That is, each frame is processed in full by EF before the next camera’s frame for that timestep is processed. Note that the ordering of the cameras here is arbitrary.

An important difference between our system and EF is that each camera has its own active and inactive map regions, $\{\Theta_i\}$ and $\{\Phi_i\}$. This allows cameras to close local loops in situations where they transition to regions of the map that are in the active regions of other cameras but that are inactive with respect to itself. To achieve this the previously scalar last seen time attribute of a surfel takes the new form of an array $t = \{t_0, \dots, t_n\}$ where t_i contains the time that the surfel was last seen by camera \mathbb{C}_i , for $i \in [1..n]$. When a measurement from frame \mathbb{F}_i^t is fused into surfel \mathcal{M}^s then the last seen time of that surfel is updated as $\mathcal{M}_i^s = t$. Each camera uses its

own active and inactive regions for mapping and tracking. For global loop closure detection, all cameras in a given *reference frame* (see below) collectively maintain a single Fern database which they use to independently check for and close global loops.

Returning now to the initial phase of processing where cameras are operating independently and the transformations of each camera relative to a shared coordinate frame are unknown. Periodically, loop closures between camera specific submaps are sought. In Section 3.1 we outline our strategy for determining inter-map loop closures. Once a loop between two submaps has been identified and the transformation between them computed the two submaps can be merged into a single map from which point multi-camera fusion can proceed as described above. We capture this situation with a *reference frame*, which is a tuple of the form $\mathcal{R}_j = \{\mathcal{M}_j, \{\mathbb{C}_i\}\}$, for surfel map \mathcal{M}_j and set of cameras $\{\mathbb{C}_i\}$ which incrementally update and track against \mathcal{M}_j through multi-camera fusion. Additionally each \mathcal{R}_j maintains its own fern database for inter and intra-map loop closures.

3.1 Fern-based Inter-Map Loop Closures

Our approach to inter-map loop closures is an extension of EF’s existing intra-map global loop closure mechanism described in Section 2. Each camera is initialised in its own reference frame. For each current active frame prediction \mathbb{F}_i^t we search the fern databases of all the other reference frames for a matching view. If one is found then we compute the transformation between the reference frame of the corresponding camera, \mathbb{C}_i , and the reference frame of the matched fern using the camera odometry from Section 2. For example, if we find a loop closure between reference frames \mathcal{R}_j and \mathcal{R}_k by matching frame \mathbb{F}_i^t , captured by \mathbb{C}_i^j in \mathcal{R}_j , to a fern in the fern encoding database of \mathcal{R}_k then the result of aligning the fern and the frame is a transformation $H_i^f \in \mathbb{SE}_3$ that aligns the matched fern and \mathbb{F}_i^t . For clarity, we have added the superscript j to the camera, which denotes the frame of reference of the camera. H_i^f gives an initial estimate of the transformation that aligns \mathcal{R}_j with \mathcal{R}_k as follows

$$\hat{T}_j^k = H_i^f (\mathbb{C}_i^j)^{-1} \quad (1)$$

To further refine the transformation and verify that it is valid we first predict a full resolution frame, ${}_k\mathbb{F}_i^t$ using \mathcal{M}_k and \mathbb{C}_i^k , which is the transformed pose of camera \mathbb{C}_i in frame of reference \mathcal{R}_k . That is, we predict the view that \mathbb{C}_i has of \mathcal{M}_k . Taking this predicted frame we align it to the original frame \mathbb{F}_k^t , again using the camera odometry described in Section 2. The result of this is a refined estimate, T_j^k , of the transformation from \mathcal{R}_j to \mathcal{R}_k . Before accepting the loop closure as a valid one we inspect the final error and inlier count of the camera tracking cost function. For the loop closure to be accepted the error must be low while the number of inliers must be high. Inliers here refer to depth measurements in \mathbb{F}_i^t which have been associated with a corresponding depth value in ${}_k\mathbb{F}_i^t$. As per the original EF algorithm the final state of the cost function covariance matrix is also checked to ensure that the cost function was sufficiently constrained along each of the 6 degrees of freedom of the camera. If all these conditions are met then the inter-map loop closure is accepted and the maps are merged.

To merge the two maps we apply T_j^k to the positions and normals of all of the surfels $\mathcal{M}^s \in \mathcal{M}_j$ and all of the camera poses, $\{\mathbb{C}_i\}$ in frame of reference \mathcal{R}_j as follows

$$\mathcal{M}_{k_p}^s = T_j^k \mathcal{M}_{j_p}^s \quad (2) \quad \mathcal{M}_{k_n}^s = R_j^k \mathcal{M}_{j_n}^s \quad (3) \quad \mathbb{C}_i^k = T_j^k \mathbb{C}_i^j \quad (4)$$

where $R_j^k \in \mathbb{SO}_3$ is formed from the upper left 3×3 submatrix of T_j^k . Note that we have not explicitly shown conversions to and from homogenous representations for the sake of brevity. This process is visualised in Figure 1. The algorithm is listed in Algorithm 1.

4 Experimental Evaluation

In benchmarking our system we aim to do two things: (i) quantitatively measure the impact of multiple sensors on mapping and tracking performance in terms of both accuracy and processing time; (ii) qualitatively

Algorithm 1 Dense collaborative mapping.

```

1: procedure COLLABORATIVEMAPPING( $c$ )
   Input: a set of  $n$  cameras  $\{C_i\}$ 
   Output: a set of  $m$  maps  $\{\mathcal{M}_j\}$  where  $m \leq n$ 

2: Dict  $c\_to\_m \leftarrow \emptyset$ 
3: for  $C_i \in c$  do
4:    $c\_to\_m[C_i] \leftarrow \text{new Map}()$ 

5: loop
6:   for  $\langle C_i, \mathcal{M}_j \rangle \in c\_to\_m$  do
7:      $\text{map\_and\_track}(C_i, \mathcal{M}_j)$ 

8:   for  $C_i \in c\_to\_m$  do
9:     Map  $\mathcal{M}_k \leftarrow c\_to\_m[C_i]$ 
10:    Mat4x4  $P_i^k \leftarrow C_i.\text{currentPose}()$ 
11:    Frame  $\mathbb{F}_{C_i^k}^P \leftarrow \text{predictView}(P_i^k, \mathcal{M}_j)$ 
12:    for  $\mathcal{M}_j \in c\_to\_m$  do
13:      if  $C_i \notin \mathcal{M}_j$  then
14:        Mat4x4  $P_i^k \leftarrow \text{findFern}(C_i, \mathcal{M}_j)$ 
15:        Frame  $\mathbb{F}_{C_i^k}^P \leftarrow \text{predictView}(P_i^k, \mathcal{M}_j)$ 
16:         $P_i^j \leftarrow \text{rgbOdometry}(\mathbb{F}_{C_i^k}^P, \mathbb{F}_{C_i^j}^P, P_i^k)$ 
17:        if  $\text{rgbOdometry}$  was successful then
18:          Mat4x4  $T_k^j \leftarrow P_i^j P_i^k^{-1}$ 
19:           $\mathcal{M}_j \leftarrow \text{merge}(\mathcal{M}_j, \mathcal{M}_k, T_k^j)$ 
20:           $c\_to\_m[C_i] \leftarrow \mathcal{M}_j$ 

```

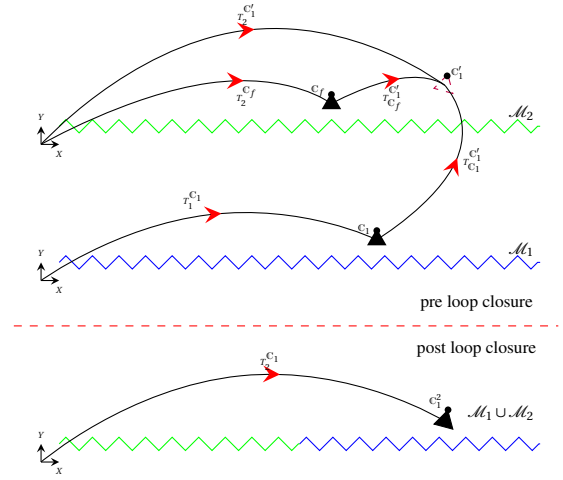


Figure 1: A fern based intermap loop closure. Camera C_1 matches its current view of \mathcal{M}_1 to a fern view, denoted by C_f , in \mathcal{M}_2 's fern database. The surfaces in each view are geometrically consistent and so a loop is closed between \mathcal{M}_1 and \mathcal{M}_2 . The loop closure itself results in $T_2^{C_1'} = T_2^{C_f} T_{C_f}^{C_1}$, the pose of C_1 in the reference frame of map \mathcal{M}_2 . Applying the composite transformation, $T_{C_1}^{C_1'} = T_2^{C_1'} (T_1^{C_1})^{-1}$, brings the pose C_1 and the surfels \mathcal{M}_1 into the coordinate frame of \mathcal{M}_2 , resulting in the fused map.

demonstrate the capability of our system. To achieve this we use two standard SLAM benchmark datasets, the ICL-NUIM dataset of [Handa et al., 2014] and the RGB-D dataset of [Sturm et al., 2012]. All experiments were run on a machine equipped with an NVidia GeForce GTX 1080 Ti GPU, 11GB of GDDR5 VRAM, a 4 core Intel i7-7700K CPU running at 4.20GHz and 16GB of DDR4 system memory. For a visualisation of dense collaborative mapping please see the accompanying video (<https://youtu.be/GUthHrKEM85M>).

We performed a quantitative analysis of the performance of the system using the living room scene from the ICL-NUIM dataset. There are four trajectories through this scene accompanied by both ground truth camera poses for each trajectory and a synthetic ground truth surface model. In our experiments, we used three of the four living room trajectories, kt_0 , kt_1 and kt_2 . We note that we have not included kt_3 in our results due to the fact that kt_3 exposes a failure mode of the system. kt_3 is one of the more challenging sequences in the ICL-NUIM dataset, as is reflected in the results of both [Dai et al., 2017] and [Whelan et al., 2016]. We found that including kt_3 in any of the combinations in our experiments resulted in significantly reduced reconstruction quality. We are currently investigating the source of this issue.

With the remaining trajectories we looked at the surface reconstruction accuracy and camera trajectory estimation accuracy of the system. To measure the surface reconstruction accuracy we processed an exhaustive set of combinations of the living room trajectories. For each combination we then computed the average per-vertex error in metres using the reference ground truth model. To measure the estimated trajectory accuracy we computed the average trajectory error root mean square error (ATE RMSE) in metres between each of the estimated trajectories in each of the combinations and the corresponding reference ground truth trajectories. Table 1 and its caption lists the specific trajectory combinations we used in these experiments and the results we recorded. Additionally we have provided pose estimation accuracy results using the TUM-RGBD dataset of [Sturm et al., 2012]. We used the ‘fr3/office’ sequence, as it is a long and challenging sequence. See Figure 2 and its caption for details of how we generated collaborative mapping sessions with this dataset. The models shown in Figure 3 demonstrate the quality of the maps produced by our system.

Following this, we again used the ‘fr3/office’ sequence, this time to measure the processing time of the

<i>Surface reconstruction and trajectory estimation accuracy for both the ICL-NUIM and TUM datasets.</i>		
Sequences	Surface Reconstruction Accuracy	Individual Trajectory ATE RMSE
ICL-NUIM		
kt_0	0.007m	0.014m
kt_1	0.008m	0.012m
kt_2	0.009m	0.027m
kt_3	0.058m	0.308m
kt_0 & kt_1	0.007m	0.016m 0.014m
kt_0 & kt_2	0.009m	0.016m 0.026m
kt_1 & kt_2	0.01m	0.0124m 0.027m
kt_0 & kt_1 & kt_2	0.009m	0.015m 0.016m 0.027m
TUM RGBD		
$office_0$	–	0.018m
$office_0$ & $office_1$	–	0.019m 0.015m
$office_0$ & $office_1$ & $office_2$	–	0.014m 0.015m 0.014m

Table 1: Surface reconstruction accuracy and trajectory estimation results from the ICL-NUIM and TUM datasets in metres¹. The first column in each row gives the specific combination of trajectories used in that session. The second column gives the per collaboration surface reconstruction accuracy. In the last column we give the per sequence ATE RMSE for each collaboration.

different components of the system during a 2 camera and then a 3 camera collaborative session. The results of this experiment are shown in Figure 2.

5 Conclusion

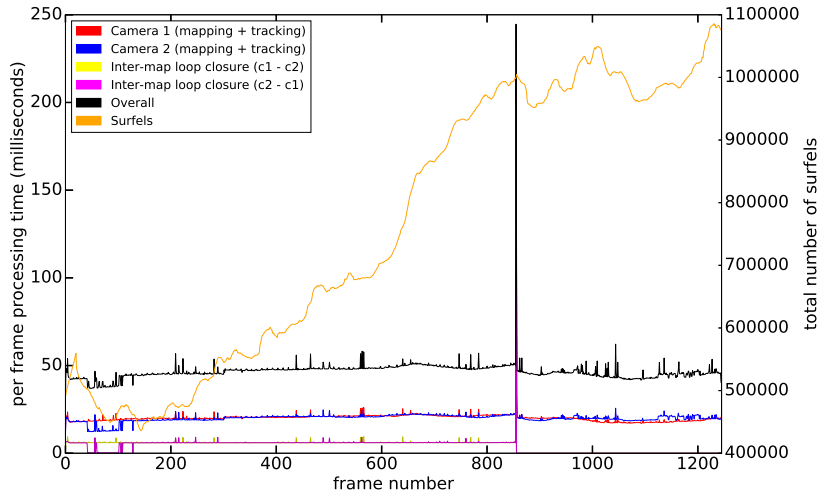
We have described a dense collaborative mapping system capable of reconstructing maps with multiple independently moving cameras, in real-time. ElasticFusion formed the foundation of our system, providing camera tracking and surface fusion capabilities. We used fern-based visual recognition to identify overlap between camera submaps and EF’s camera tracking mechanism to find the transformations between the coordinate frames of the matched submaps.

We provided a quantitative analysis of our system using the ICL-NUIM dataset to measure its surface reconstruction accuracy and the camera trajectory estimation accuracy. We then used the Freiburg dataset to give a breakdown of frame processing times for the system.

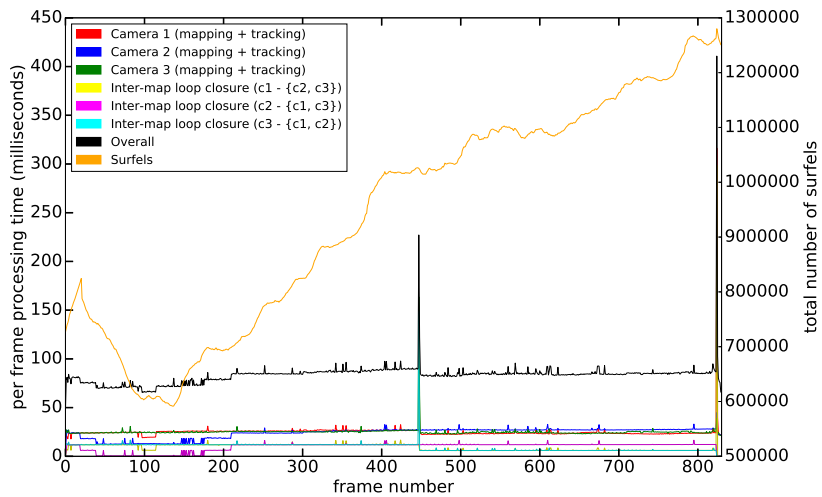
Currently the system works well for mapping with 2 or 3 cameras, however, mapping with a larger number of cameras would involve overcoming significant challenges. For example, to maintain real-time operation, a more sophisticated approach to mapping and tracking, involving processing each frame in parallel, would be required. Another direction for future work is the modelling of scene dynamics, a pertinent problem in the context of collaborative mapping as cameras that are reconstructing a scene would inevitably move into and out of view of one another, almost guaranteeing some level of scene dynamics.

Recall from Section 2 that the deformation graph influence region for a surfel is computed based on a *temporal-spatial* neighborhood. However, in a collaborative session, cameras which are close in time may be fusing measurements into the model that are far apart in space, meaning that searching the graph for temporally nearby nodes may not yield a spatially nearby influence region with the knock on effect of inhibiting deformation graph optimisation. Thus the influence region computation needs to be rethought in future work.

¹The reader will note that in the single camera sessions our results are not exactly as reported in [Whelan et al., 2016]. This is for a number of reasons; (i) the system is sensitive to the specific GPU of the machine it is deployed on; (ii) as fern encoding is a random process loop closures will not always occur at exactly the same point.



(a)



(b)

Figure 2: Timings for (a) a 2 camera and (b) a 3 camera collaborative mapping session with the ‘fr3/office’ sequence. The original sequence contained 2487 frames. For graph (a) we split these frames into two subsequences of equal length, similarly, for graph (b) we split them into three subsequences of equal length. The subsequences were then processed collaboratively. For example, in (a) Camera 1 corresponds to the first half of the frames from the original sequence while Camera 2 corresponds to the second half of the frames. The spike at frame 942 in graph (a) is caused by an inter-map loop closure between the two cameras. Similarly the two spikes in graph (b) at frames 449 and 826 are both caused by inter-map loop closures.

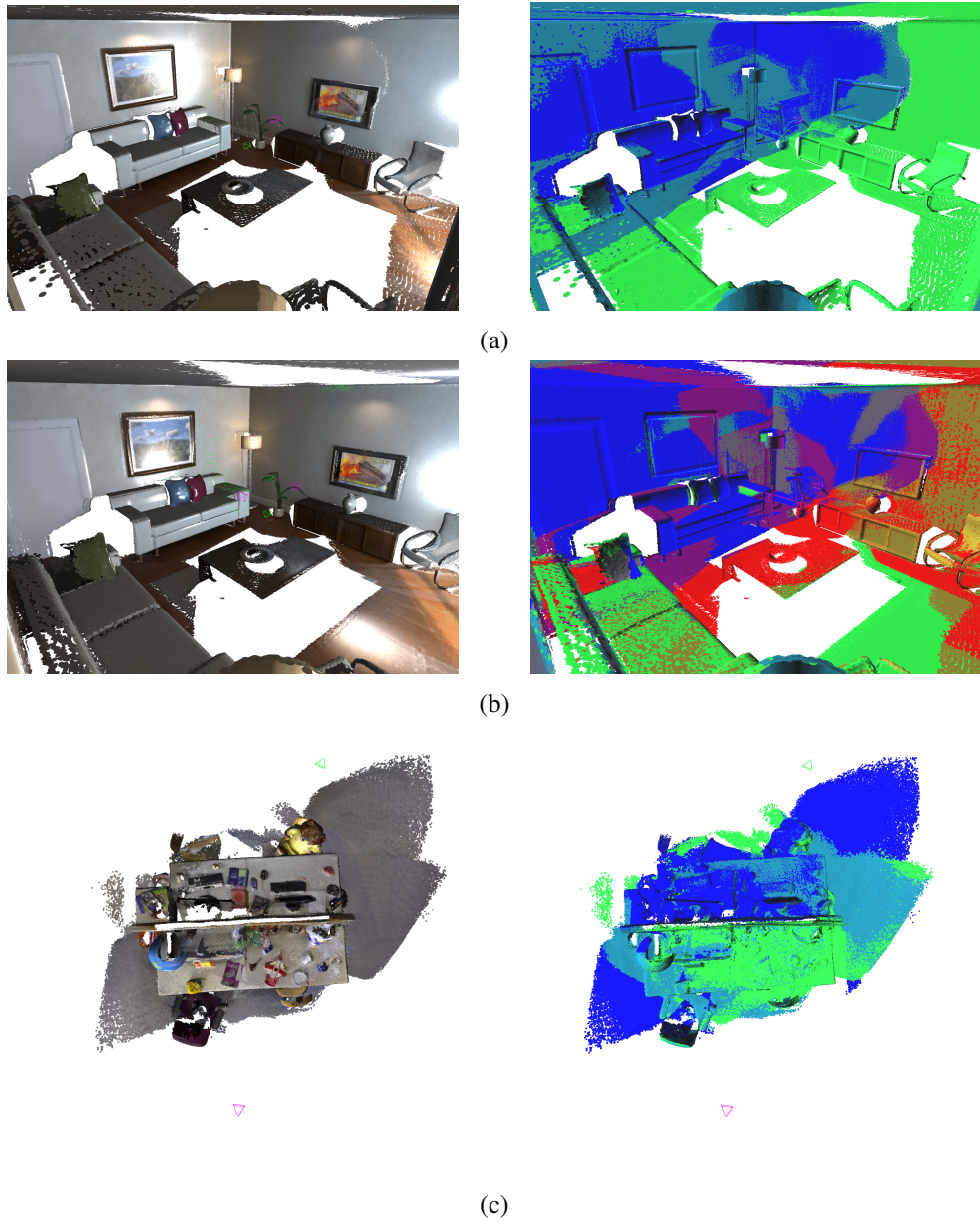


Figure 3: Examples of maps produced by the system. The left column shows the final models produced by the system. The right column shows the contributions of each camera to the final model. Each cameras specific contributions are coloured with a separate primary colour. Composite coloured surfels have fused measurements from multiple cameras. Rows (a) and (b) correspond to the ICL-NUIM kt_0 & Kt_2 and kt_0 & kt_1 & kt_2 sessions respectively. Row (c) corresponds to the 2 camera session using the ‘fr3/office’ sequence from the Freiburg dataset.

References

- [Dai et al., 2017] Dai, A., Nießner, M., Zollhöfer, M., Izadi, S., and Theobalt, C. (2017). Bundlefusion: Real-time globally consistent 3D reconstruction using on-the-fly surface reintegration. *ACM Trans. Graph.*, 36(3):24:1–24:18.
- [Engel et al., 2014] Engel, J., Schöps, T., and Cremers, D. (2014). LSD-SLAM: Large-scale direct monocular SLAM. In *European Conference on Computer Vision (ECCV)*.
- [Fenwick et al., 2002] Fenwick, J., Newman, P., and Leonard, J. J. (2002). Cooperative concurrent mapping and localization. In *Proceedings of the 2002 IEEE International Conference on Robotics and Automation, ICRA 2002, May 11-15, 2002, Washington, DC, USA*.
- [Glocker et al., 2015] Glocker, B., Shotton, J., Criminisi, A., and Izadi, S. (2015). Real-time RGB-D camera relocalization via randomized ferns for keyframe encoding. *IEEE Transactions on Visualization and Computer Graphics*, 21(5):571–583.
- [Golodetz et al., 2018] Golodetz, S., Cavallari, T., Lord, N. A., Prisacariu, V. A., Murray, D. W., and Torr, P. H. S. (2018). Collaborative large-scale dense 3D reconstruction with online inter-agent pose optimisation. *CoRR*, abs/1801.08361.
- [Handa et al., 2014] Handa, A., Whelan, T., McDonald, J., and Davison, A. (2014). A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM. In *IEEE Intl. Conf. on Robotics and Automation, ICRA, Hong Kong, China*.
- [Howard, 2006] Howard, A. (2006). Multi-robot simultaneous localization and mapping using particle filters. *Int. J. Rob. Res.*, 25(12):1243–1256.
- [Howard et al., 2006] Howard, A., Sukhatme, G. S., and Mataric, M. J. (2006). Multirobot simultaneous localization and mapping using manifold representations. *Proceedings of the IEEE*, 94(7):1360–1369.
- [Huang et al., 2010] Huang, A. S., Olson, E., and Moore, D. C. (2010). LCM: Lightweight communications and marshalling. In *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 4057–4062.
- [Keller et al., 2013] Keller, M., Lefloch, D., Lambers, M., Izadi, S., Weyrich, T., and Kolb, A. (2013). Real-time 3D reconstruction in dynamic scenes using point-based fusion. In *Proceedings of the 2013 International Conference on 3D Vision, 3DV '13*, pages 1–8, Washington, DC, USA. IEEE Computer Society.
- [Kerl et al., 2013] Kerl, C., Sturm, J., and Cremers, D. (2013). Dense visual SLAM for RGB-D cameras. In *Proc. of the Int. Conf. on Intelligent Robot Systems (IROS)*.
- [McDonald et al., 2013] McDonald, J., Kaess, M., Cadena, C., Neira, J., and Leonard, J. (2013). Real-time 6-DOF multi-session visual SLAM over large-scale environments. *Robotics and Autonomous Systems*, 61(10):1145–1158.
- [Mur-Artal and Tardos, 2017] Mur-Artal, R. and Tardos, J. (2017). ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. 33(5):1255–1262.
- [Newcombe et al., 2011a] Newcombe, R. A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A. J., Kohli, P., Shotton, J., Hodges, S., and Fitzgibbon, A. (2011a). Kinectfusion: Real-time dense surface mapping and tracking. In *Proceedings of the 2011 10th IEEE International Symposium on Mixed and Augmented Reality, ISMAR '11*, pages 127–136, Washington, DC, USA. IEEE Computer Society.

- [Newcombe et al., 2011b] Newcombe, R. A., Lovegrove, S. J., and Davison, A. J. (2011b). DTAM: Dense tracking and mapping in real-time. In *Proceedings of the 2011 International Conference on Computer Vision, ICCV '11*, pages 2320–2327, Washington, DC, USA. IEEE Computer Society.
- [Prisacariu et al., 2017] Prisacariu, V. A., Kähler, O., Golodetz, S., Sapienza, M., Cavallari, T., Torr, P. H., and Murray, D. W. (2017). InfiniTAM v3: A Framework for Large-Scale 3D Reconstruction with Loop Closure. *ArXiv e-prints*.
- [Saeedi et al., 2016] Saeedi, S., Trentini, M., Seto, M., and Li, H. (2016). Multiple-robot simultaneous localization and mapping: A review. *J. Field Robot.*, 33(1):3–46.
- [Sturm et al., 2012] Sturm, J., Engelhard, N., Endres, F., Burgard, W., and Cremers, D. (2012). A benchmark for the evaluation of RGB-D SLAM systems. In *Proc. of the International Conference on Intelligent Robot Systems (IROS)*.
- [Thrun and Liu, 2003] Thrun, S. and Liu, Y. (2003). Multi-robot slam with sparse extended information filters. In *Springer Tracts in Advanced Robotics*, volume 15, pages 254–266.
- [Whelan et al., 2014] Whelan, T., Kaess, M., Johannsson, H., Fallon, M., Leonard, J., and McDonald, J. (2014). Real-time large scale dense RGB-D SLAM with volumetric fusion. *Intl. J. of Robotics Research, IJRR*.
- [Whelan et al., 2016] Whelan, T., Salas-Moreno, R. F., Glocker, B., Davison, A. J., and Leutenegger, S. (2016). Elasticfusion: Real-time dense SLAM and light source estimation. *Intl. J. of Robotics Research, IJRR*.
- [Williams et al., 2002] Williams, S. B., Dissanayake, G., and Durrant-Whyte, H. (2002). Towards multi-vehicle simultaneous localisation and mapping. In *Proceedings 2002 IEEE International Conference on Robotics and Automation*, volume 3, pages 2743–2748.
- [Zhou and Roumeliotis, 2006] Zhou, X. S. and Roumeliotis, S. I. (2006). Multi-robot SLAM with unknown initial correspondence: The robot rendezvous case. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1785–1792.