

# Real-time large scale dense RGB-D SLAM with volumetric fusion

Thomas Whelan, Michael Kaess, Hordur Johannsson, Maurice Fallon, John J. Leonard and John McDonald

## Abstract

We present a new SLAM system capable of producing high quality globally consistent surface reconstructions over hundreds of metres in real-time with only a low-cost commodity RGB-D sensor. By using a fused volumetric surface reconstruction we achieve a much higher quality map over what would be achieved using raw RGB-D point clouds. In this paper we highlight three key techniques associated with applying a volumetric fusion-based mapping system to the SLAM problem in real-time. First, the use of a GPU-based 3D cyclical buffer trick to efficiently extend dense every frame volumetric fusion of depth maps to function over an unbounded spatial region. Second, overcoming camera pose estimation limitations in a wide variety of environments by combining both dense geometric and photometric camera pose constraints. Third, efficiently updating the dense map according to place recognition and subsequent loop closure constraints by the use of an “as-rigid-as-possible” space deformation. We present results on a wide variety of aspects of the system and show through evaluation on *de facto* standard RGB-D benchmarks that our system performs strongly in terms of trajectory estimation, map quality and computational performance in comparison to other state-of-the-art systems.

**Keywords:** volumetric fusion, camera pose estimation, dense methods, large scale, real-time, RGB-D, SLAM, GPU

## 1 Introduction

The ability for a robot to create a map of an unknown environment and localise within that map is of extreme importance in intelligent autonomous operation. Simultaneous Localisation and Mapping (SLAM) has been one of the large focuses of robotics research over the last two decades, with 3D mapping becoming more and more popular within the last few years over traditional 2D laser scan SLAM. The recent explosion in full dense 3D SLAM is arguably a result of the release of the Microsoft Kinect commodity RGB-D sensor, which provides high quality depth sensing capabilities for a little over one hundred US dollars. Before the advent of the Kinect, 3D SLAM methods required either time of flight (TOF) sensors, 3D LIDAR scanners or stereo vision, which were typically either quite expensive or not suitable for fully mobile real-time operation if dense reconstruction was desired. Another recent technology which is often coupled with dense methods is General-Purpose computing on Graphics Processing Units

(GPGPU) which exploits the massive parallelism available in GPU hardware to perform high speed and often real-time processing on entire images every frame. Being an affordable commodity technology, GPU-based programming is arguably another large enabler in recent dense SLAM research.

Many visual SLAM systems and 3D reconstruction systems (both offline and online) have been published in recent times that rely purely on RGB-D sensing capabilities because of the Kinect’s low price and accuracy; Henry et al. (2012); Endres et al. (2012); Stückler and Behnke (2013). The Kinect-Fusion algorithm of Newcombe et al. (2011) is one of the most notable RGB-D-based 3D reconstruction systems of recent times, allowing real-time volumetric dense reconstruction of a desk sized scene at sub-centimetre resolution. By fusing many individual depth maps together into a single volumetric reconstruction, the models that are obtained are of much higher quality than typical noisy single-shot raw RGB-D point clouds. KinectFusion enables reconstructions of an unprecedented quality at real-time speeds but comes with a number of limitations, namely 1) restriction to a fixed small area in space; 2) reliance on geometric information alone for camera pose estimation; and, 3) no means of explicitly incorporating loop closures. These three limitations severely limit the applicability of KinectFusion to the large scale SLAM problem where it is desirable due to its real-time nature and very high surface reconstruction fidelity.

In this paper we present solutions to the three aforementioned limitations such that the system can be used in a full real-time large scale SLAM setting. We address the three limitations respectively by 1) representing the volumetric re-

---

T. Whelan and J. McDonald are with the Department of Computer Science, National University of Ireland Maynooth, Co. Kildare, Ireland. [thomas.j.whelan@nuim.ie](mailto:thomas.j.whelan@nuim.ie), [johnmcd@cs.nuim.ie](mailto:johnmcd@cs.nuim.ie)

M. Kaess is with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA. [kaess@cmu.edu](mailto:kaess@cmu.edu)

H. Johannsson, M. Fallon and J. Leonard are with the Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts 02139, USA. [{hordurj,mfallon,jleonard}@mit.edu](mailto:{hordurj,mfallon,jleonard}@mit.edu)

This work was presented in part at the Robotics Science and Systems RGB-D Workshop, Sydney, July 2012 (Whelan et al. (2012)), in part at the International Conference on Robotics and Automation, Karlsruhe, May 2013 (Whelan et al. (2013a)) and in part at the International Conference on Intelligent Robots and Systems, Japan, November 2013 (Whelan et al. (2013b)).

construction data structure in memory with a rolling cyclical buffer; 2) estimating a dense photometric camera constraint in conjunction with a dense geometric constraint and jointly optimising for a camera pose estimate; and, 3) optimising the dense map by means of a non-rigid space deformation parameterised by a loop closure constraint. In the remainder of this section we provide a discussion on the existing work related to the area of dense RGB-D SLAM. Following on from this Sections 2, 3 & 4 address the issues of extended scale volumetric fusion, camera pose estimate, and loop closure, respectively. Section 5 provides a comprehensive qualitative and quantitative evaluation of the system using multiple benchmark datasets and finally Section 6 presents conclusions on the work and future directions of our research.

## 1.1 Related Work

A large number of publications have been made over the last few years specifically using RGB-D data for camera pose estimation, dense mapping and full SLAM pipelines. While many visual SLAM systems existed prior to the advent of active RGB-D sensors (e.g. Comport et al. (2007)), we will focus mainly on the literature which makes specific use of active RGB-D platforms. One of the earliest RGB-D tracking and mapping systems uses FAST feature correspondences between frames for visual odometry and offloads dense point cloud map building to a post-processing step utilising sparse bundle adjustment (SBA) for global consistency by minimizing feature reprojection error (Huang et al. (2011)). One of the first real-time dense RGB-D tracking and mapping systems estimates an image warping function with both geometric and photometric information to compute a camera pose estimate, however only relies on rigid reprojection for point cloud map reconstruction without using a method for global consistency (Audras et al. (2011)). Similar work on dense RGB-D camera tracking was done by Steinbrücker et al. (2011), also estimating an image warping function based on geometric and photometric information. Recent work by Kerl et al. (2013) presents a more robust dense photometrics-based RGB-D visual odometry system that proposes a t-distribution-based error model which more accurately matches the residual error between RGB-D frames in scenes that are not entirely static.

Henry et al. (2012) presented one of the first full SLAM systems based entirely upon RGB-D data, using visual feature matching with Generalised Iterative Closest Point (GICP) to build up a pose graph and following that an optimised surfel map of the area explored. The use of pose graph optimisation versus SBA is studied, minimising feature reprojection error in an offline rigid transformation framework. Visual feature correspondences are used in conjunction with pose graph optimisation in the RGB-D SLAM system of Endres et al. (2012). An octree-based volumetric representation is used to store the map, created by reprojecting all point measurements into the global frame. This map representation is provided by the OctoMap framework of Hornung et al. (2013), which includes the ability to take measurement uncertainties into ac-

count and implicitly represent free and occupied space while being space efficient. An explicit voxel volumetric occupancy representation is used by Pirker et al. (2011) in their GPSlam system which uses sparse visual feature correspondences for camera pose estimation. They make use of visual place recognition and sliding window bundle adjustment in a pose graph optimisation framework. To achieve global consistency the occupancy grid is “morphed” by a weighted average of the log-odds perceptions of each camera for each voxel. Stückler and Behnke (2013) register surfel maps together for camera pose estimation and store a multi-resolution surfel map in an octree, using pose graph optimisation for global consistency. After pose graph optimisation is complete a globally consistent map is created by fusing key views together. In recent work Hu et al. (2012) proposed a system that uses bundle adjustment in order to make use of pixels for which no valid depth exists, and Lee et al. (2012) presented a system which exploits GPU processing power for real-time camera tracking. Both systems produce an optimised map as a final step in the process.

A substantial number of derived works have been published recently after the advent of the KinectFusion system of Newcombe et al. (2011), mostly focused on extending the range of operation, with other related work on object recognition and motion planning (Karpthy et al. (2013); Wagner et al. (2013)). Recent work by Bylow et al. (2013) and Canelhas et al. (2013) directly tracks the camera pose against the accumulated volumetric model by exploiting the fact that the truncated signed distance function (TSDF) representation used by KinectFusion stores the signed distance to the closest surface at voxels near the surface. This avoids the need to raycast a vertex map for each frame to perform camera pose estimation, which potentially discards information about the surface reconstruction.

Roth and Vona (2012) extend the operational range of KinectFusion by using a double buffering mechanism to map between volumetric models upon camera translation and rotation, using a voxel interpolation for the latter. However no method for recovering the map is provided. Zeng et al. (2012) replace the explicit voxel representation used by KinectFusion with an octree representation which allows mapping of areas up to 8m×8m×8m in size. However this method does increase the chance for drift within the map and provides no means of loop closure or map correction. Steinbrücker et al. (2013) make use of a multi-scale octree to represent the signed distance function, allowing full color reconstructions of scenes as large as an entire corridor including nine rooms spanning a total area of 45m×12m×3.4m. After an RGB-D sequence has been processed, a globally consistent camera trajectory is resolved and the model is reconstructed. Keller et al. (2013) present an extended fusion system made space efficient by using a point-based surfel representation, although lacking in drift correction or loop closure detection. Chen et al. (2013) present a novel hierarchical data structure that enables extremely space efficient volumetric fusion, using a streaming

framework allowing effectively unbounded mapping range, limited only by available memory. However the system lacks any method for mitigating drift or enforcing global consistency. Nießner et al. (2013) present an alternative space efficient method for large scale dense fusion that uses an intelligent voxel hashing function to minimise the amount of memory required for reconstruction, but again without a means of correcting for drift.

An alternative approach to the modern SLAM problem is introduced by Salas-Moreno et al. (2013), whereby known objects are detected, tracked and mapped in real-time in a dense RGB-D framework. Pose graph optimisation is used to ensure global consistency on the level of camera poses and detected object positions. This does allow loop closure, however less influence is placed on a full scene reconstruction with only point cloud reprojections being used for mapped loop closure. Recent work by Henry et al. (2013b) uses multiple smaller “patch volumes” to segment the mapped space into a set of discrete TSDFs, each with a 6-degrees-of-freedom (6-DOF) pose which is rigidly optimised upon loop closure detection. This approach can be seen as similar to the SLAM++ approach of Salas-Moreno et al. (2013) whereby the patch volumes are analogous to objects. While achieving global consistency between each volume, there is no clear solution presented for correcting the surface within any one given volume or stitching surfaces which are split between volumes, leaving local surfaces disconnected.

Zhou et al. (2013) present an impressive method for reconstructing 3D scenes that specifically targets the high-frequency noise and low-frequency distortion effects often encountered with RGB-D data. By reconstructing fragments of the scene which are then aligned and deformed very high quality reconstructions can be obtained, however in what is a strictly offline framework. Similar work also by Zhou and Koltun (2013) presents a method which detects points of interest in a scene and specifically optimises the camera trajectory to preserve detailed geometry around these points, within an offline frame.

An number of approaches that rely on keyframes have been developed to tackle the problem of RGB-D mapping and SLAM. Tykkälä et al. (2013) present a system which uses real-time dense photometric keyframe-based camera tracking to determine a camera trajectory around an indoor environment. Individual RGB-D frames are also fused into existing keyframes to improve reconstruction quality. An optional bundle adjustment step can then be taken to optimise the camera poses before a watertight Poisson mesh reconstruction is computed as a post-processing step. Meilland and Comport (2013) propose a model that unifies the benefits of a dense voxel-based representation with a keyframe representation allowing high quality dense mapping over large-scales, although without detecting large loop closures or correcting for drift. An intelligent forward composition approach is proposed which enables efficient combination of reference images to create a single predicted frame without repeated

redundant image warps. In our work we chose to avoid a keyframe approach in spite of the resulting higher memory requirement. A fully 3D voxel-based method removes the need to implement specific schemes to overcome the problems associated with reconstructing complex non-concave objects and non-convex scenes.

As discussed there exists a large number of systems utilising RGB-D data for SLAM and related problems. However, most are either unable to operate in real-time, provide an up-to-date optimised representation of the map at runtime or any time it is requested or efficiently incorporate large non-rigid updates to the map. Non-rigid surface correction is of great interest specifically in the realm of volumetric fusion as typically reconstructions are locally highly accurate but drift slowly over large scales over time, where a smooth continuous deformation of the surface is most suitable for correction. In the following sections we will fully describe our approach to RGB-D SLAM with volumetric fusion which is capable of functioning in real-time over large scale trajectories, while efficiently applying non-rigid updates to the dense map upon loop closure to ensure global consistency.

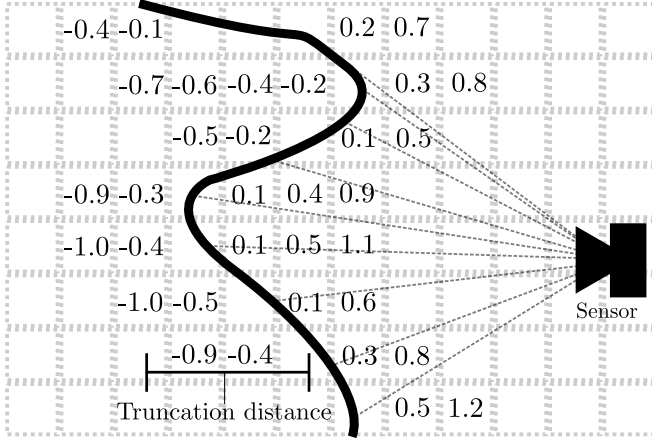
To clarify our definition of “real-time” there is no offline step involved in our pipeline and multiple loops can be closed immediately as they occur during the mapping process (shown in Multimedia Extension 2). Our system architecture can be compared to that of PTAM (Klein and Murray (2007)), whereby camera tracking and mapping run in separate threads. While the camera tracking component runs at frame rate in one thread, the mapping component is freed from the computational burden of updating the map for every frame and instead occasionally receives information from the tracking thread to update the map for consistency.

This paper brings together work presented in our three previous publications Whelan et al. (2012), Whelan et al. (2013a) and Whelan et al. (2013b). In this paper we provide a number of additions to that work including a method for improving camera-frustum overlap for greater reconstruction range (Section 2.4) and a means of reducing the amount of information required to perform map deformation, increasing computational performance (Section 5.3.2). Most significantly this paper presents an extensive evaluation of the presented system not present in any previous work, including both qualitative and quantitative evaluation of trajectory estimation performance, surface reconstruction quality and computational performance.

Please note any provided sample parameter and threshold values are those which were used for all experiments in this paper and are provided as an aid to those who wish to re-implement any aspect of this work.

## 2 Extended Scale Volumetric Fusion

In this section we will provide some background on the usage of volumetric fusion for dense RGB-D-based tracking and mapping and describe our extension to KinectFusion, the

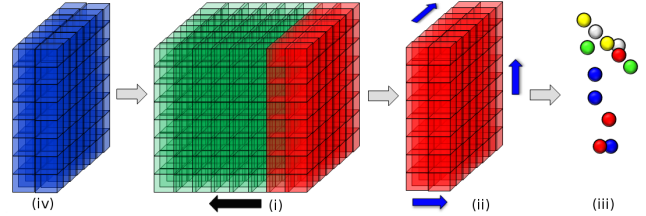


**Figure 1:** Two dimensional example of the structure of the truncated signed distance function representation of an implicit surface. Shown are example signed distance values stored at voxels within the truncation distance of the observed surface, with rays cast from the observing sensor.

most widely cited system that employs this approach, to allow spatially extended mapping.

## 2.1 Background

Real-time volumetric fusion with RGB-D cameras was brought to the forefront by Newcombe et al. (2011) with the KinectFusion system. A significant component of the system is the cyclical pipeline used for camera tracking and scene mapping, whereby full depth maps are fused into a volumetric data structure (TSDF), which is then raycast to produce a predicted surface that the subsequently captured depth map is matched against using ICP. The truncated signed distance function (TSDF) is a volumetric data structure that encodes implicit surfaces by storing the signed distance to the closest surface at each voxel up to a given truncation distance from the actual surface position. Points at which the sign of the distance value changes are known as zero crossings, which represent the actual position of the surface, shown in Figure 1. Each voxel also stores a weight for the distance measurement at that point, effectively providing a moving average of the surface position. In the case of KinectFusion, the TSDF is stored as a three dimensional voxel grid in GPU memory where dense depth map integration is accomplished by sweeping through the volume and updating distance measurements accordingly, while surface raycasting is carried out by simply projecting rays from the current camera pose and returning the depth and surface normals at the first zero crossings encountered. Surface normals are easily computed by taking the finite difference around a given position within the TSDF, as exploited by Bylow et al. (2013) and Canelhas et al. (2013). The entire process is very amenable to parallelisation and greatly benefits in execution time from being implemented on a GPU (Newcombe et al. (2011)). A point to note is that the TSDF representation has a minimal surface thickness limitation imposed by the selected truncation distance. This problem was



**Figure 2:** Visualisation of the volume shifting process for spatially extended mapping; (i) The camera motion exceeds the movement threshold  $m_s$  (direction of camera motion shown by the black arrow); (ii) Volume slice leaving the volume (red) is raycast along all three axes to extract surface points and reset to free space; (iii) The raycast surface is extracted as a point cloud and fed into the Greedy Projection Triangulation (GPT) algorithm of Marton et al. (2009); (iv) New region of space (blue) enters the volume and is integrated using new modulo addressing of the volume.

highlighted and explored by Henry et al. (2013a) in their work on multiple fusion volumes.

## 2.2 Volume Representation

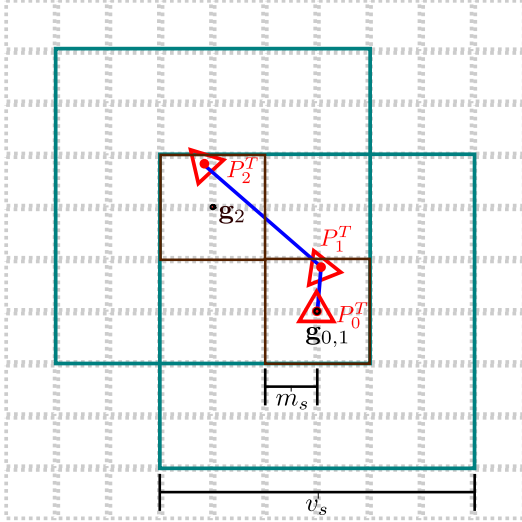
Defining the voxel space domain as  $\Psi \subset \mathbb{N}^3$  the TSDF volume  $S$  at some location  $\mathbf{s} \in \Psi$  has the mapping  $S(\mathbf{s}) : \Psi \rightarrow \mathbb{R} \times \mathbb{N} \times \mathbb{N}^3$ . Within GPU memory the TSDF is represented as a 3D array of voxels. Each voxel contains a signed distance value ( $S(\mathbf{s})^T$ , truncated float16), an unsigned weight value ( $S(\mathbf{s})^W$ , unsigned int8) and a byte for each color component R, G and B ( $S(\mathbf{s})^R, S(\mathbf{s})^G, S(\mathbf{s})^B$ ) for a total of 6 bytes per voxel. The integration of new surface measurements is carried out in a similar fashion to Newcombe et al. (2011), when integrating a new signed distance function measurement  $S(\mathbf{s})_i^T$  during the fusion of a new depth map, each voxel  $\mathbf{s} \in \Psi$  at time  $i$  is updated with:

$$S(\mathbf{s})_i^T = \frac{S(\mathbf{s})_{i-1}^W S(\mathbf{s})_{i-1}^T + S(\mathbf{s})_i^W S(\mathbf{s})_i^T}{S(\mathbf{s})_{i-1}^W + S(\mathbf{s})_i^W} \quad (1)$$

$$S(\mathbf{s})_i^W = \min(S(\mathbf{s})_{i-1}^W + S(\mathbf{s})_i^W, max\_weight) \quad (2)$$

As is the case with previous approaches, we take  $S(\mathbf{s})_i^W = 1$  to provide a simple moving average, and set  $max\_weight$  to 128. Bylow et al. (2013) have experimented with different weighting schemes, however we have found the original value of 1 used by Newcombe et al. (2011) to provide good performance. Using only a cubic volume, we parameterise the TSDF by the side length in voxels  $v_s$  and the dimension in metres  $v_d$ . Both of these parameters control the resolution of the reconstruction along with the size of the immediate “active area” of reconstruction. In all experiments in this paper we set  $v_s = 512$  for total GPU memory usage of 768MB. The 6-DOF camera pose within the TSDF at time  $i$  is denoted as  $P_i^T$ , composed of a rotation  $\mathbf{R}_i^T \in \mathbb{SO}_3$  and a translation  $\mathbf{t}_i^T \in \mathbb{R}^3$ . The origin of the TSDF coordinate system is positioned at the center of the volume with basis vectors aligned with the axes of the TSDF. Initially  $\mathbf{R}_0^T = \mathbf{I}$  and  $\mathbf{t}_0^T = (0, 0, 0)^T$ . The position of the TSDF volume in voxel units in the global frame is ini-





**Figure 3:** Visualisation of the interaction between the movement threshold  $m_s$  and the shifting process. Between frames 0 and 1 the camera does not cross the movement boundary (dark brown) and no shift occurs. At frame 2, the pose crosses the boundary and causes a volume shift, recentering the volume (teal) around  $P_2^T$  and updating  $\mathbf{g}_2$ . The underlying voxel grid quantisation is shown in light dashed lines.

tialised to be  $\mathbf{g}_0 = (0, 0, 0)^T$ . Note that the superscript  $T$  refers to the TSDF pose and not the transpose  $\top$  operator.

### 2.3 Volume Shifting

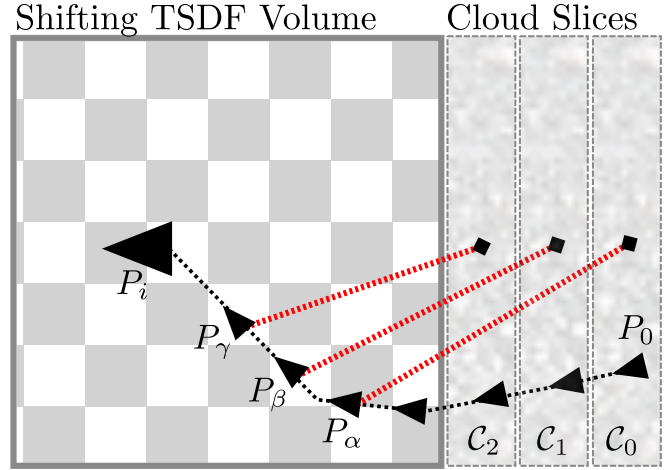
Unlike Newcombe et al. (2011) camera pose estimation and surface reconstruction is not restricted to only the region around which the TSDF was initialised. By employing modulo arithmetic in how the TSDF volume is addressed in GPU memory we can treat the structure like a cyclical buffer which virtually translates as the camera moves through an environment. Figure 2 provides a visual example and description of the shifting process. It is parameterised by an integer movement threshold  $m_s$ , defining the cubic movement boundary (in voxels) around  $\mathbf{g}_i$  which upon crossing, causes a volume shift, shown in Figure 3. Discussion on the choice of value for  $m_s$  is provided in Section 5.3. Each dimension is treated independently during a shift. When a shift is triggered, the TSDF is virtually translated about the camera pose (in voxel units) to bring the camera’s position to within one voxel of  $\mathbf{g}_{i+1}$ . The new pose of the camera  $P_{i+1}^T$  has no change in rotation, while the shift corrected camera position  $\mathbf{t}_{i+1}^{T'}$  is calculated from  $\mathbf{t}_{i+1}^T$  by first computing the number of voxel units crossed:

$$\mathbf{u} = \left\lfloor \frac{v_s \mathbf{t}_{i+1}^T}{v_d} \right\rfloor \quad (3)$$

And then shifting the pose while updating the global position of the TSDF:

$$\mathbf{t}_{i+1}^{T'} = \mathbf{t}_{i+1}^T - \frac{v_d \mathbf{u}}{v_s} \quad (4)$$

$$\mathbf{g}_{i+1} = \mathbf{g}_i + \mathbf{u} \quad (5)$$



**Figure 4:** Two dimensional visualisation of the association between extracted cloud slices, the camera poses and the TSDF volume. Note that the camera poses here are in global coordinates rather than internal TSDF coordinates. A red dashed line links camera poses with extracted slices of the TSDF volume ( $P_\gamma, P_\beta$  and  $P_\alpha$  with  $C_2, C_1$  and  $C_0$  respectively). The large triangles represent camera poses that caused volume shifts while the small black squares represent those that didn’t.

#### 2.3.1 Implementation

There are two parts of volumetric fusion as described by Newcombe et al. (2011) that require indexed access to the TSDF volume; 1) Volume Integration and 2) Volume Raycasting. Referring again to Figure 2, the new surface measurements shown in blue can be integrated into the memory previously used for the old surface contained within the red region of the TSDF by ensuring all element look ups in the 3D block of GPU memory reflect the virtual voxel translation computed in Equation 5. Assuming row major memory ordering, an element in the unshifted cubic 3D voxel grid can be found at the 1D memory location  $a$  given by:

$$a = (x + yv_s + zv_s^2) \quad (6)$$

The volume’s translation can be reflected in how the TSDF is addressed for integration and raycasting by substituting the indices in Equation 6 with values that are offset by the current global position of the TSDF and bound within the dimensions of the voxel grid using the modulus operator:

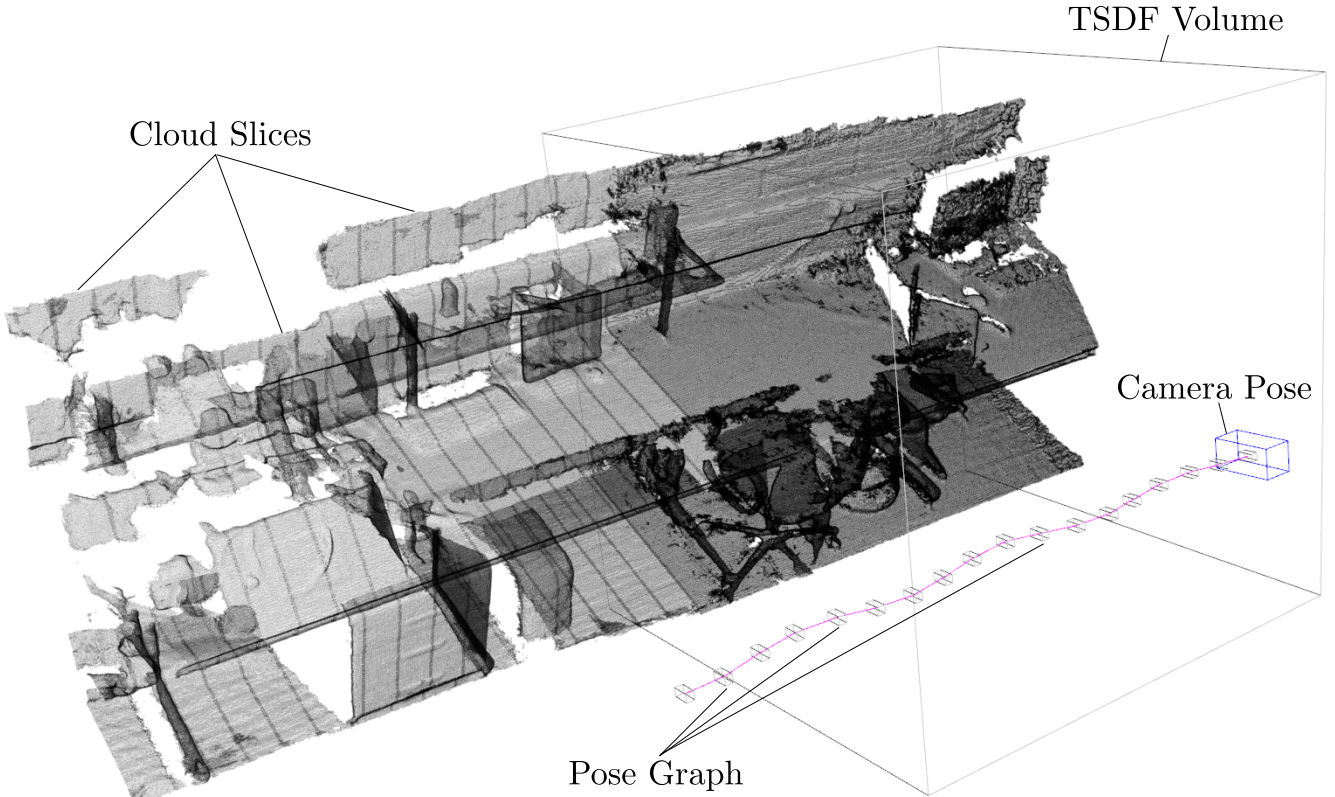
$$x' = (x + \mathbf{g}_{ix}) \mod v_s \quad (7)$$

$$y' = (y + \mathbf{g}_{iy}) \mod v_s \quad (8)$$

$$z' = (z + \mathbf{g}_{iz}) \mod v_s \quad (9)$$

$$a = (x' + y'v_s + z'v_s^2) \quad (10)$$

The original KinectFusion work of Newcombe et al. (2011) benefits greatly from memory caching and pipelining functionality within GPU memory to achieve high computational performance within the integration step. In our implementation we have found that use of a cyclical addressing method



**Figure 5:** Visualisation of a shifted TSDF volume with extracted cloud slices and pose graph highlighted, using dynamic cube positioning discussed in Section 2.4. The pose graph is drawn in pink, while small cuboids are drawn for camera poses that have cloud slices associated with them. Note that the apparent striping of the boundaries between slices has been added in for visualisation purposes only.

has no significant effect on real-time performance. An explanation for the lack of a drastic performance decrease is that even after significant buffer cycling there still exists continuous blocks of memory which at least partially maintain pipelining.

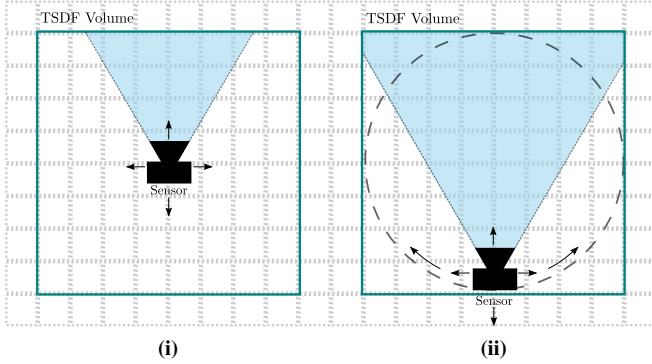
### 2.3.2 Surface Extraction

In order to recover the surface from the TSDF that moves out of the region of space encompassed by the volume the  $\mathbf{u}$  value computed in Equation 3 is used with  $\mathbf{g}_i$  to index a three dimensional slice of the volume to extract surface points from. These points are extracted by three orthogonal raycasts aligned with the axes of the TSDF through the slice, extracting zero crossings as individual surface vertices. We filter out noisy measurements at this point by only extracting points that have a minimum voxel weight. The same 3D slice of the volume is then reset to free space to allow integration of new surface measurements. The extracted vertices are transferred to main system memory where further processing takes place. The orthogonal raycast can result in duplicate vertices if the TSDF is obliquely aligned to the surface being reconstructed. A voxel grid filter is used to remove these points by overlaying a voxel grid (with the same voxel size as the TSDF) on the extracted point cloud and returning a new point cloud with a point for each voxel that represents the centroid

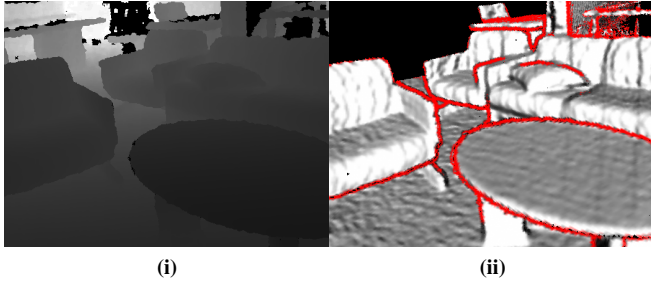
of all points that fell inside that voxel. Each set of vertices extracted from the TSDF in this fashion is known as a “cloud slice”. From here, we rebuild the surface by incrementally triangulating successive cloud slices using an incremental mesh growing variant of the GPT algorithm to ensure surface connectivity between slices (Marton et al. (2009)).

We choose not to perform marching cubes because this would lock the TSDF data structure in GPU memory and delay the reset of the extracted volume slice, impacting volume shifting performance over all. Axis-aligned orthogonal raycasting is extremely fast and allows us to offload the data from the GPU and unlock the TSDF volume as quickly as possible. This way the GPU-based tracking and integration components of the system can continue with minimal interruption while the extracted cloud slice is triangulated on the CPU asynchronously. In addition to this the raycast vertex map representation is easy to incrementally grow using the mesh triangulation method we have adopted (Marton et al. (2009)).

We associate with each cloud slice the pose of the camera at the time of the slice’s extraction. This is visualised in Figure 4. At this point we introduce camera poses in the global coordinate frame outside of the TSDF volume  $P_i$ , composed of a rotation  $\mathbf{R}_i \in \mathbb{S}\mathbb{O}_3$  and a translation  $\mathbf{t}_i \in \mathbb{R}^3$ . The global pose  $P_i$  of a camera from the TSDF at time  $i$  is made up of:



**Figure 6:** Visualisation of frustum-volume overlap for regular and dynamic cube positioning, from left to right; (i) By keeping the camera centered in the TSDF, there is poor overlap between the camera’s field of view and the volume; (ii) By using a circular (or spherical) parameterisation of the volume’s position relative to the camera, greater overlap with and usage of the TSDF can be achieved.



**Figure 7:** From left to right; (i) Input depth map registered to RGB channel; (ii) Color measurements from pixels highlighted in red are rejected for being on depth discontinuities. Lighter surfaces are weighted higher in color integration due to being well aligned with the camera sensor.

$$\mathbf{R}_i = \mathbf{R}_i^T \quad (11)$$

$$\mathbf{t}_i = \mathbf{t}_i^T + \frac{v_d \mathbf{g}_i}{v_s} \quad (12)$$

We construct a pose graph incrementally using each global camera pose  $P_i$ , that is, a camera pose for every frame where some poses are attached to cloud slices. The full shifting and surface extraction process is shown in Figure 5, where only the poses with associated cloud slices are drawn.

## 2.4 Dynamic Cube Positioning

As mentioned in Section 2.2, we position the camera in the center of the TSDF volume and roughly maintain this position inside the TSDF at all times. This parameterisation of the camera position relative to the volume is wasteful as most of the volume is unused (i.e. behind the camera) and there is little overlap between the camera frustum and the volume, shown in Figure 6. By dynamically changing the position of the volume relative to the camera depending on the camera’s orientation we can achieve greater frustum-volume overlap

and make better use of the entire TSDF volume. In a typical SLAM setting a circular parameterisation is sufficient.

Defining  $\beta_i$  to be the rotation around the  $y$ -axis of the camera pose at time  $i$ , we can compute the new position of the center of the TSDF volume relative to the camera as:

$$\mathbf{r}^T = \left( \frac{v_d}{2} \cdot \cos\left(\beta_i + \frac{\pi}{2}\right), 0, \frac{v_d}{2} \cdot \sin\left(\beta_i - \frac{\pi}{2}\right) \right)^T \quad (13)$$

This dynamic parameterisation enables more intelligent use of the volume and maintains a larger active reconstruction area in front of the camera at all times, while also being easily expandable to a full spherical parameterisation depending on the expected camera motion.

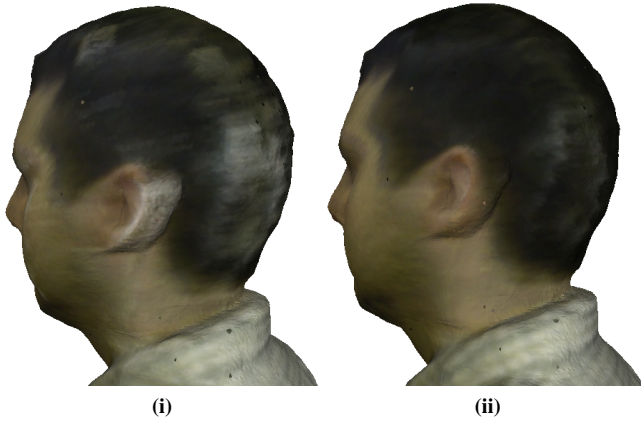
## 2.5 Color Estimation

As well as estimating the surface itself in the reconstruction process, we also estimate the color of the surface (purely for visualisation purposes). Color is integrated into the TSDF in a similar manner to depth measurements including value truncation and averaging. The only distinction is that the predicted surface color values obtained from the volume raycast are not used in camera pose estimation. The motivation for this decision is discussed into Section 3.2. Color fusion has similar advantages to depth map fusion in that sensor noise and other optical phenomena are averaged out from the final reconstruction over time.

### 2.5.1 Artifact Reduction

The estimated surface color is sometimes inaccurate around the edges of closed objects in a scene due to poor calibration between the RGB and depth cameras or light diffraction around objects. We have observed that there typically exists stark discontinuities in the depth channel around such edges which can in turn cause the background to blend with the foreground surface or vice-versa. To address this issue we opt to reject the integration of color measurements close to or on strong boundaries in the depth image. A color measurement is deemed to be on a boundary if some of its neighbours are more than a given distance away from it in depth. We consider a pixel neighbourhood window of  $7 \times 7$  pixels around each RGB value to be integrated. Figure 7 shows a source depth image and rejected measurements on the TSDF surface model. In addition to this it is ideal to weight color measurements on surfaces well aligned with the sensor higher than those at extreme angles. We weight each color measurement update by the normal angle on the surface with respect to the sensor, visualised in Figure 7. The more parallel the surface is to the image plane, the higher the weight on the color measurement.

Defining the image space domain as  $\Omega \subset \mathbb{N}^2$ , an RGB-D frame  $I_i$  is composed of an RGB image  $\mathbf{rgb}_i : \Omega \rightarrow \mathbb{N}^3$ , a depth image  $\mathbf{d}_i : \Omega \rightarrow \mathbb{R}$  and a timestamp  $i$ . We also define a normal map computed for  $\mathbf{d}_i$  as  $\mathbf{n}_i : \Omega \rightarrow \mathbb{R}^3$ . We list the



**Figure 8:** From left to right; (i) Light diffraction behind a foreground surface has caused incorrect color integration (ii) Incorporating a discontinuity check with surface angle weighting greatly reduces the visual artifacts captured.

algorithm for color integration in Algorithm 1 in Appendix B. Note that we define the  $z$ -axis to point outward from the sensor and in all experiments use an RGB-D frame resolution of  $640 \times 480$ . An example reconstruction is shown in Figure 8 comparing surface coloring with and without the described measures.

### 3 Camera Pose Estimation

A number of volumetric fusion systems use only depth information for camera pose estimation (Newcombe et al. (2011), Chen et al. (2013), Bylow et al. (2013), Keller et al. (2013), Roth and Vona (2012), Zeng et al. (2012), Canelhas et al. (2013)). A reliance on geometric information alone for camera pose estimation has a number of well understood problems, such as the inability to function in corridor-like environments and other scenes with few 3D features. To avoid these problems like Henry et al. (2013b) we combine dense geometric camera pose constraints with dense photometric constraints to achieve a more robust pose estimate in more challenging scenes. We base our approach on the dense photometric image warping method of Steinbrücker et al. (2011) and Audras et al. (2011), performing dense RGB-D alignment every frame in real-time. In tune with other components of the pipeline we utilise a GPU implementation of the algorithm. Following we describe the geometric and photometric components of the camera pose estimation pipeline and our method for combining them to form a single joint pose constraint.

#### 3.1 Geometric Camera Pose Estimation

Many of the previous works on volumetric fusion estimate the pose of the camera each frame relative to the TSDF by aligning the current depth map with the TSDF, either by ray-casting the volume to retrieve a vertex and normal map of the predicted surface (as done originally by Newcombe et al.

(2011)) and performing iterative closest point (ICP) or by directly minimising the distance to the surface in the TSDF (Bylow et al. (2013), Canelhas et al. (2013)). We perform the former in order to avoid expensive global memory accesses in the TSDF volume in GPU memory.

We aim to find the motion parameters  $\xi$  that minimize the cost over the point-to-plane error between vertices in the current depth frame and the predicted raycast surface:

$$E_{icp} = \sum_k \left\| \left( \mathbf{v}^k - \exp(\hat{\xi}) \mathbf{T} \mathbf{v}_n^k \right) \cdot \mathbf{n}^k \right\|^2, \quad (14)$$

where  $\mathbf{v}_n^k$  is the  $k$ -th vertex in frame  $n$ ,  $\mathbf{v}^k, \mathbf{n}^k$  are the corresponding vertex and normal in the model, and  $\mathbf{T}$  is the current estimate of the transformation from the current frame to the model frame. For simplicity of notation we omit conversions between 3-vectors (as needed for dot and cross products) and their corresponding homogeneous 4-vectors (as needed for multiplications with  $\mathbf{T}$ ). We utilise projective data association as originally proposed by Newcombe et al. (2011) for fast point correspondence between the vertex maps by projecting the vertices from the predicted surface  $\mathbf{v}$  onto the depth map vertices  $\mathbf{v}_n$ . Linearizing the transformation around the identity we get:

$$E_{icp} \approx \sum_k \left\| \left( \mathbf{v}^k - (\mathbf{I} + \hat{\xi}) \mathbf{T} \mathbf{v}_n^k \right) \cdot \mathbf{n}^k \right\|^2 \quad (15)$$

$$= \sum_k \left\| \left( \mathbf{v}^k - \mathbf{T} \mathbf{v}_n^k \right) \cdot \mathbf{n}^k - \hat{\xi} \mathbf{T} \mathbf{v}_n^k \cdot \mathbf{n}^k \right\|^2 \quad (16)$$

$$= \sum_k \left\| \begin{bmatrix} -\mathbf{T} \mathbf{v}_n^k \times \mathbf{n}^k \\ -\mathbf{n}^k \end{bmatrix}^T \xi + (\mathbf{v}^k - \mathbf{v}_n^k) \cdot \mathbf{n}^k \right\|^2 \quad (17)$$

$$= \left\| \mathbf{J}_{icp} \xi + \mathbf{r}_{icp} \right\|^2 \quad (18)$$

Blocks of the measurement Jacobian and residual can be populated in tandem and solved with a highly parallel tree reduction on the GPU to produce a  $6 \times 6$  system of normal equations which are then transferred to the CPU and solved with Cholesky decomposition to yield  $\hat{\xi}$ . As in previous work we compute the alignment iteratively with a three level coarse-to-fine depth map pyramid scheme.

#### 3.2 Photometric Camera Pose Estimation

As mentioned previously we choose to match between consecutive RGB-D frames with the photometric component instead of matching to the predict surface reconstruction. Depending on the configuration of the TSDF there may be poor overlap between the camera frustum and the volume, which limits the amount of photometric information which can be used, where distant photometric features are desirable to constrain camera rotation. As well as this, the resolution of the TSDF in terms of voxels may produce a ray-cast image with a much lower resolution than the image produced by the RGB sensor. By default the Microsoft Kinect



and Asus Xtion Pro Live, two of the most popular RGB-D sensors, have automatic exposure and white balance enabled, which can cause unusual coloring of the surface reconstruction over time, again hindering model-based photometric tracking. While these functions of the camera can be disabled we have found that it is sometimes desirable to keep them enabled, particularly in indoor environments where lighting can vary to a certain degree.

Given two consecutive RGB-D frames  $[\mathbf{rgb}_{n-1}, \mathbf{d}_{n-1}]$  and  $[\mathbf{rgb}_n, \mathbf{d}_n]$  we compute a rigid camera transformation between the two that maximises photoconsistency. Defining  $V : \Omega \rightarrow \mathbb{R}^3$  to be the back-projection of a point  $\mathbf{p}$ , dependent on a metric depth map  $M : \Omega \rightarrow \mathbb{R}$  and camera intrinsics matrix  $\mathbf{K}$  made up of the principal points  $c_x$  and  $c_y$  and the focal lengths  $f_x$  and  $f_y$ :

$$V(\mathbf{p}) = \left( \frac{(\mathbf{p}_x - c_x)M(\mathbf{p})}{f_x}, \frac{(\mathbf{p}_y - c_y)M(\mathbf{p})}{f_y}, M(\mathbf{p}) \right)^\top \quad (19)$$

We also defined perspective projection of a 3D point  $\mathbf{v} = (x, y, z)^\top$  including dehomogenisation by  $\Pi(\mathbf{v}) = (x/z, y/z)^\top$ . The cost we wish to minimise depends on the difference in intensity values between two images  $I_{n-1}, I_n : \Omega \rightarrow \mathbb{N}$ :

$$E_{rgb} = \sum_{\mathbf{p} \in \mathcal{L}} \left\| I_n(\mathbf{p}) - I_{n-1}(\Pi_{n-1}(\exp(\hat{\xi})\mathbf{T}V_n(\mathbf{p}))) \right\|^2 \quad (20)$$

Where  $\mathcal{L}$  is the list of valid interest points populated in Algorithm 2 (see Appendix B) and  $\mathbf{T}$  is the current estimate of the transformation from  $I_n$  to  $I_{n-1}$ . Similar to the geometric pose estimation method we solve for this transformation iteratively with a three level image pyramid.

### 3.2.1 Preprocessing

For both pairs we perform preprocessing on the RGB image and depth map. For each depth map we convert raw sensor values to a metric depth map  $M : \Omega \rightarrow \mathbb{R}$  and we compute an intensity image  $I = (\mathbf{rgb}^R * 0.299 + \mathbf{rgb}^G * 0.587 + \mathbf{rgb}^B * 0.114)$  with  $I : \Omega \rightarrow \mathbb{N}$ . Following this a three level intensity and depth pyramid is constructed using a  $5 \times 5$  Gaussian kernel for downsampling. We compute the partial derivatives  $\frac{\partial I_n}{\partial x}$  and  $\frac{\partial I_n}{\partial y}$  using a  $3 \times 3$  Sobel operator coupled with a  $3 \times 3$  Gaussian blur with  $\sigma = 0.8$ . Each of these steps is carried out on the GPU acting in parallel with one GPU thread per pixel.

### 3.2.2 Precomputation

As with the ICP method described in Section 3.1, we use projective data association between frames to population point correspondences. For the sake of speed we only include point correspondences with a minimum gradient in the intensity image, with the motivation that other low gradient points will not have a significant effect on the final transformation. We implement this optimisation by using a list of interest points,

which involves a much larger set of points than a point feature extractor could provide. Compiling this list of points as a parallel operation is done using a basic parallel reduction exploiting shared memory in each CUDA thread block as inspired by a similar operation by van den Braak et al. (2011). Algorithm 2 in Appendix B lists the operation as it would operate for each level of the pyramid.

In the computation of the Jacobian matrix the projection of each point in  $M_{n-1}$  is required. For each pyramid level the 3D projection  $V_{n-1}(\mathbf{p})$  of each point  $\mathbf{p}$  in the depth map is computed prior to beginning iteration. Only projecting certain points based on a condition results in performance hindering branching and a reduction in pipelining. Empirically it was found to be faster to simply project the entire depth map rather than only project points required in correspondences.

### 3.2.3 Iterative Transformation Estimation

Our iterative estimation process takes two main steps; (i) populating a list of valid correspondences from the precomputed list of interest points and (ii) solving the linear system for an incremental transformation and concatenating these transformations. The first step involves a reduction similar to the one in Algorithm 2, but rather than reducing from a 2D array to a 1D array it reduces from a 1D array to another 1D array; a distinction which results in a notable difference in implementation. On the first iteration for frame  $n$  we set the estimated camera transformation matrix  $\mathbf{T}$  to the identity, where

$$\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 0 & 0 & 1 \end{bmatrix} \in \mathbb{SE}_3 \quad (21)$$

with a rotation  $\mathbf{R} \in \mathbb{SO}_3$  and translation  $\mathbf{t} \in \mathbb{R}^3$ . Before each iteration we compute the projection of  $\mathbf{T}$  into the image before uploading to the GPU as

$$\mathbf{R}' = \mathbf{K}\mathbf{R}\mathbf{K}^{-1}, \quad \mathbf{t}' = \mathbf{K}\mathbf{t}. \quad (22)$$

Algorithm 3 in Appendix B lists the process of populating a list of point correspondences from the list of interest points which can then be used to construct the Jacobian. With a list of valid correspondences we need only solve a least-squares equation

$$\arg \min_{\xi} \left\| \mathbf{J}_{rgb} \xi + \mathbf{r}_{rgb} \right\|^2 \quad (23)$$

to compute an improved camera transformation estimate

$$\mathbf{T}' = \exp(\hat{\xi})\mathbf{T} \quad (24)$$

$$\hat{\xi} = \begin{bmatrix} [\omega]_{\times} & \mathbf{x} \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (25)$$

with  $\xi = [\omega^\top \mathbf{x}^\top]^\top$ ,  $\omega \in \mathbb{R}^3$  and  $\mathbf{x} \in \mathbb{R}^3$ . We first normalise the intensity difference sum  $\sigma$  computed in Algorithm 3 to enable a weighted optimisation  $\sigma' = \sqrt{\sigma/k_C}$ . Computation of the  $\sigma$  value in parallel is in fact an optimisation exploiting the atomic arithmetic functions available in the CUDA API. From here  $\mathbf{J}_{rgb}$  and  $\mathbf{r}_{rgb}$  can be populated including usage

of  $\sigma'$  for weighting. Equation 23 is then solved using a tree reduction on the GPU followed by Cholesky factorisation of the linear system on the CPU.

### 3.3 Combined Camera Pose Estimate

We combine the cost functions of both the geometric and photometric estimates in a weighted sum, independent of the number of points used for each estimate. The sum of the RGB-D and ICP cost is defined as

$$E = E_{icp} + w_{rgb}E_{rgb} \quad (26)$$

where  $w_{rgb}$  is the weight and was set empirically to 0.1 to reflect the difference in metrics used for ICP and RGB-D costs (metres and 8-bit intensity respectively). A key distinction between our approach and that of Tykkälä et al. (2011) is that we are combining two cost functions between a frame-to-model registration (for the ICP component) and a frame-to-frame registration (for the RGB-D component). For each step we minimize the linear least-squares problem by solving the normal equations

$$\begin{bmatrix} \mathbf{J}_{icp} \\ v\mathbf{J}_{rgb} \end{bmatrix}^T \begin{bmatrix} \mathbf{J}_{icp} \\ v\mathbf{J}_{rgb} \end{bmatrix} \xi = \begin{bmatrix} \mathbf{J}_{icp} \\ v\mathbf{J}_{rgb} \end{bmatrix}^T \begin{bmatrix} \mathbf{r}_{icp} \\ \mathbf{r}_{rgb} \end{bmatrix} \quad (27)$$

$$(\mathbf{J}_{icp}^T \mathbf{J}_{icp} + w_{rgb} \mathbf{J}_{rgb}^T \mathbf{J}_{rgb}) \xi = \mathbf{J}_{icp}^T \mathbf{r}_{icp} + v \mathbf{J}_{rgb}^T \mathbf{r}_{rgb} \quad (28)$$

where  $v = \sqrt{w_{rgb}}$ . The products  $\mathbf{J}^T \mathbf{J}$  and  $\mathbf{J}^T \mathbf{r}$  are computed on the GPU using a tree reduction. The normal equations are then solved on the CPU using Cholesky factorisation. The final estimate returns a locally optimal (in the least-squares sense) camera pose which jointly minimizes the photometric error between the current RGB-D frame and the last and the geometric error between the current depth map and the TSDF surface reconstruction. This combined method provides a very accurate and stable trajectory estimate as well as surface reconstruction, which we expand upon in Section 5.

It should be noted that although there are a number of atomic operations in Algorithms 2 and 3, these are primarily operating on values contained in shared thread block memory, minimising impact on execution performance and overall degradation to serial execution. Our computational performance results in Section 5.3 and our previous work (Whelan et al. (2013a)) also demonstrates that use of such atomic operations (in the standard reduction setting they are used in here) does not hinder real-time performance.

## 4 Loop Closure

Using the techniques from Sections 2 and 3 permits the reconstruction of large scale dense 3D mesh-based maps in real-time, however like all egomotion estimation systems drift will accumulate over space and time, warranting a method to correct the map to achieve global consistency when possible. A simple approach to this problem would be to associate each vertex in the mesh with the nearest camera pose, optimise the

pose graph and reflect the camera pose transformations in the mesh vertices. This would however cause sharp discontinuities at points on the surface where the association between camera poses changes and ignores other important properties of the surface. For this reason we have chosen a non-rigid method of correcting the map. We now frame the system as a more traditional SLAM setup with a frontend (for camera tracking and surface extraction) and a backend (for pose graph optimisation and map optimisation). A detailed system architecture diagram is shown in Figure 9.

The frontend is made up of the extended scale volumetric fusion method described in Section 2 coupled with the combined geometric and photometric camera pose estimation method described in Section 3. The final component of the frontend not yet described is a visual place recognition module that relies on the DBoW place recognition system (Galvez-Lopez and Tardos (2011)) which we describe in Section 4.2.

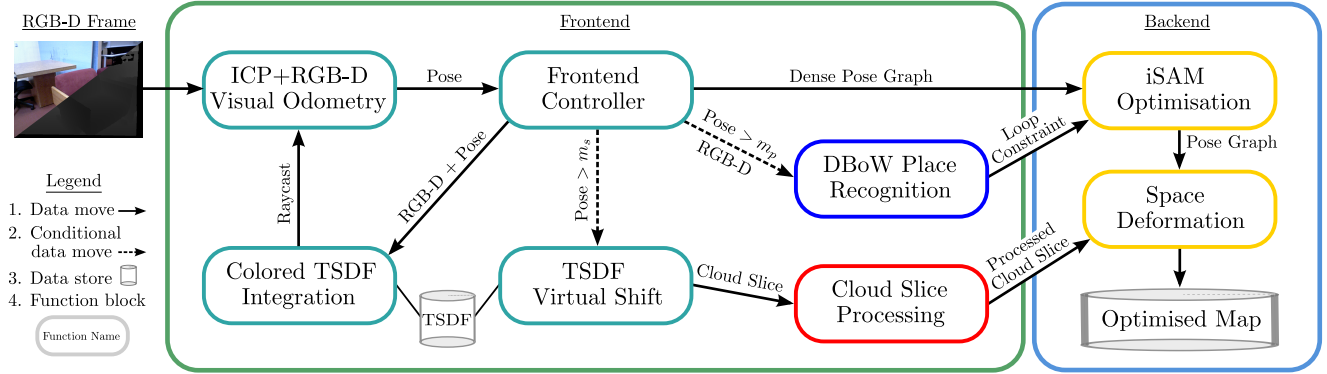
The backend provides a means of performing deformation-based dense map correction making use of incremental pose graph optimisation coupled with a non-rigid map optimisation. We use iSAM (Kaess et al. (2008)) to optimise the camera pose graph according to loop closure constraints provided by our place recognition module. The optimised trajectory is then used in conjunction with matched visual features to constrain a non-rigid space deformation of the map. We adapt the embedded deformation technique of Sumner et al. (2007) to apply it to large scale dense maps captured with a pose graph backend and utilise efficient incremental methods to prepare the map for deformation.

We apply the SLAM principal to our framework by building constraints between multiple regions of the surface through frames anchored to the map via the place recognition system. These frames (and associated global camera poses) are connected to the pose graph, which upon optimisation propagates back to the surface through the deformation.

Following we provide a detailed description of each component involved in the global consistency pipeline including pose graph representation, place recognition and loop closure, deformation graph construction and map optimisation.

### 4.1 Pose Graph

All camera poses added to the pose graph are given in global coordinates, as described in Section 2.3.2. A camera pose  $P_i$  is estimated for every processed frame. We evaluate the trade offs of using every pose versus a subset of poses in Section 5. As discussed in Section 2.3.2 some camera poses also have an associated cloud slice as shown in Figure 10 where the relationship between pose  $P_\gamma$  and cloud slice  $C_j$  is shown. This provides a useful association between camera poses and the extracted surface, capturing both temporal and spatial proximity. In order to model the uncertainty of inter-pose constraints derived from dense visual odometry we can approximate the constraint uncertainty with the Hessian as  $\Sigma = (\mathbf{J}^T \mathbf{J})^{-1}$ , where



**Figure 9:** System architecture diagram. Differently colored function blocks are executing asynchronously in separate CPU threads. The  $m_s$  quantity denotes the volume shifting threshold and  $m_p$  denotes the place recognition movement threshold.

$\mathbf{J}$  is the combined measurement Jacobian computed in Equation 28.

## 4.2 Place Recognition

We use Speeded Up Robust Feature (SURF) descriptors with the bag-of-words-based DBoW loop detector for place recognition (Galvez-Lopez and Tardos (2011)). Adding every RGB-D frame to the place recognition system is non-optimal, therefore we utilise a movement metric sensitive to both rotation and translation which indicates when to add a new frame to the place recognition system. Defining  $\mathbf{r}(\mathbf{R}) : \mathbb{SO}(3) \rightarrow \mathbb{R}^3$  to provide the rotation vector form of some rotation matrix  $\mathbf{R}$ , we compute a movement distance between two poses  $a$  and  $b$  that compounds both translation and rotation into a single quantity as:

$$m_{ab} = \|\mathbf{r}(\mathbf{R}_a^{-1}\mathbf{R}_b)\|_2 + \|\mathbf{t}_a - \mathbf{t}_b\|_2 \quad (29)$$

For each frame we evaluate the movement distance between the current frame pose and the pose of the last frame added to the place recognition system according to Equation 29. If this metric is above some threshold  $m_p$ , a new frame is added. Empirically we found  $m_p = 0.3$  provides good performance. Alternatively the two quantities can be separately thresholded such that motion is acknowledged when either  $\|\mathbf{r}(\mathbf{R}_a^{-1}\mathbf{R}_b)\|_2$  goes above a specified angle  $\theta_t$  threshold or  $\|\mathbf{t}_a - \mathbf{t}_b\|_2$  goes above a distance  $m_t$  threshold. We have not found place recognition rates to vary significantly between schemes.

Upon receiving a new RGB-D frame  $[\mathbf{rgb}_i, \mathbf{d}_i]$  the place recognition module first computes a set of SURF keypoints and associated descriptors  $U_i \in \Omega \times \mathbb{R}^{64}$  for that frame. These features are cached in memory for future queries. The depth image  $\mathbf{d}_i$  is also cached, however to ensure low memory usage it is compressed in real-time using lossless compression (Deutsch and Gailly (1996)). Following this, the existing bag-of-words descriptor database is queried. If a match is found the SURF keypoints and descriptors  $U_m$  and depth data  $\mathbf{d}_m$  for the matched image are retrieved for constraint computa-

tion. A number of validation steps are performed to minimise the chance of false positives. They are as follows:

### 4.2.1 SURF Correspondence Threshold

Given  $U_i$  and  $U_m$  we find correspondences by a k-nearest neighbour search in the SURF descriptor space. We use the Fast Library for Approximate Nearest Neighbors (FLANN) to perform this search and populate a set of valid correspondences  $G \in \Omega \times \Omega$ , thresholding matches using an  $L_2$ -norm between descriptors in  $\mathbb{R}^{64}$ . We discard the loop closure candidate if  $|G|$  is less than some threshold; a value of 35 has been found to provide adequate performance in our experiments.

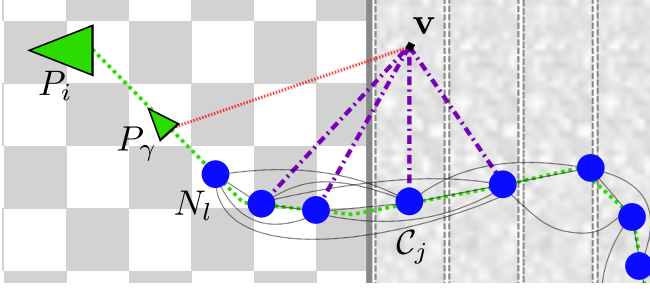
### 4.2.2 RANSAC Transformation Estimation

Given  $G$  and  $\mathbf{d}_m$ , we first attempt to approximate a 6-DOF relative transformation between the camera poses of frames  $i$  and  $m$  using a RANSAC-based 3-point algorithm (Fischler and Bolles (1981)). Each matching keypoint in  $G$  is back-projected from image  $m$  to a 3D point, transformed according to the current RANSAC model and reprojected into the image plane of frame  $i$  (using standard perspective projection onto an image plane) where the reprojection error quantified by the  $L_2$ -norm in  $\mathbb{R}^2$  is used for outlier detection. Empirically we chose a maximum reprojection error of 2.0 pixels for inliers. If the percentage of inliers for the RANSAC estimation is below 25% the loop closure is discarded. Otherwise, we refine the estimated transformation by minimising all inlier feature reprojection errors in a Levenberg-Marquardt optimisation.

### 4.2.3 Point Cloud ICP

At this point only candidate loop closures with strong geometrically consistent visual feature correspondences remain. As a final step we perform a non-linear ICP step between  $\mathbf{d}_i$  and  $\mathbf{d}_m$ . Firstly we back-project each point in both depth images to produce two point clouds. In order to speed up the computation, we carry out a uniform downsampling of each point cloud in  $\mathbb{R}^3$  using a voxel grid filter. Finally,





**Figure 10:** Two-dimensional example showing the current position of the TSDF shifting volume as a checkerboard pattern and the previously extracted cloud slices as textured columns. Also shown is the pose graph as small green points as well as a pose  $P_\gamma$  which caused a volume shift. The association between  $P_\gamma$  and the extracted cloud slice is shown with a dotted red line. A  $k = 4$  connected sequential deformation graph is also shown, demonstrating the back-traversal vertex association algorithm on a random vertex  $v$ .

using the RANSAC approximate transformation estimate as an initial guess, we iteratively minimise nearest neighbour correspondence distances between the two point clouds using a Levenberg-Marquardt optimisation. We accept the final refined transformation if the mean  $L_2^2$ -norm of all correspondence errors is below a threshold. Typically we found a threshold of 0.01 to provide good results.

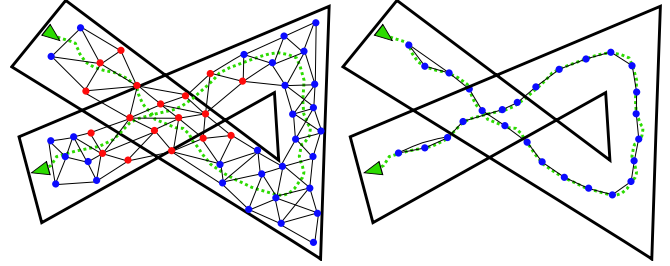
Once a loop closure candidate has passed all of the described tests, the relative transformation constraint between the two camera poses is added to the pose graph maintained by the iSAM module. Section 4.4 describes how this constraint is used to update the map.

### 4.3 Space Deformation

Our approach to non-rigid space deformation of the map is based on the embedded deformation approach of Sumner et al. (2007). Their system allows deformation of open triangular meshes and point clouds; no connectivity information is required as is the case with many deformation algorithms (Karan (2000); Jacobson and Sorkine (2011)). Exploiting this characteristic, Chen et al. (2012) applied embedded deformation to automatic skeletonised rigging and real-time animation of arbitrary objects in their KinÊtre system. Next we describe our adaptation of embedded deformation to apply it to large scale dense maps with a focus on automatic incremental deformation graph construction.

#### 4.3.1 Deformation Graph

Sumner et al. (2007) propose the use of a deformation graph to facilitate space deformation of a set of vertices. A deformation graph is composed of nodes and edges spread across the surface to be deformed. Each node  $N_l$  has an associated position  $N_l^g \in \mathbb{R}^3$  and set of neighbouring nodes  $\mathcal{N}(N_l)$ . The neighbours of each node are what make up the edges of the graph. Each node also stores an affine transformation in the form of a  $3 \times 3$  matrix  $N_l^R$  and a  $3 \times 1$  vector  $N_l^t$ , initialised by



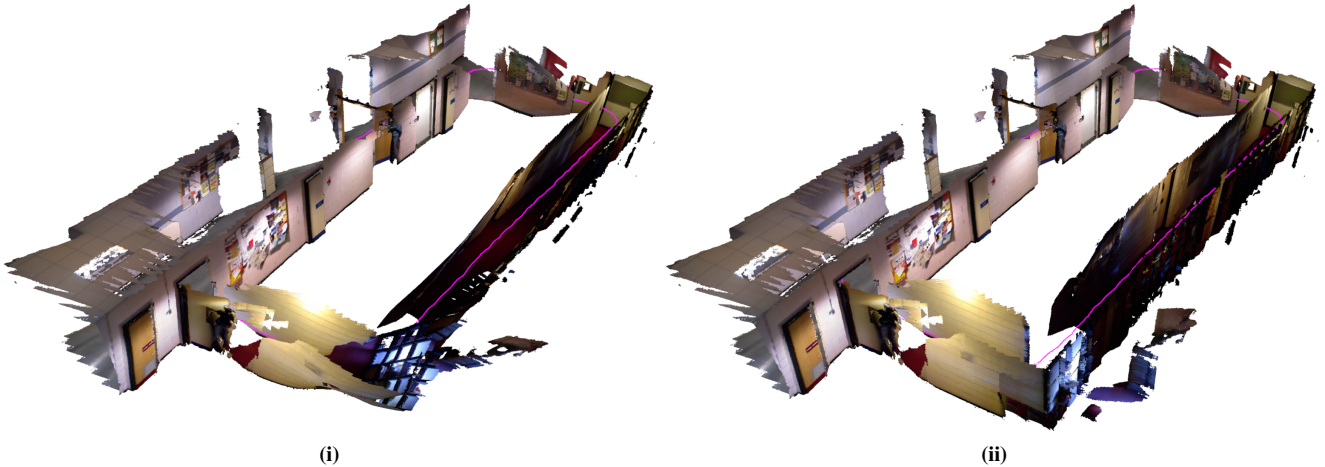
**Figure 11:** Two-dimensional example of deformation graph construction. On the left a spatially-constrained graph is constructed over a pre-loop closure map suffering from significant drift. The nodes highlighted in red are connected to nodes which belong in potentially unrelated areas of the map. On the right our incremental sampling and connectivity strategy is shown (two-nearest neighbours for simplicity) which samples and connects nodes along the pose graph, preventing unrelated areas of the map being connected by the deformation graph.

default to the identity and  $(0, 0, 0)^T$  respectively. The effect of this affine transformation on any vertex which that node influences is centered at the node’s position  $N_l^g$ .

#### 4.3.2 Incremental Graph Construction

The original approach to embedded deformation relies on a uniform sampling of the vertices in  $\mathbb{R}^3$  to construct the deformation graph. Chen et al. (2012) substitute this with a method that uses a 5D orientation-aware sampling strategy based on the Mahalanobis distance between surface points in order to prevent links in the graph between physically unrelated areas of the model. Neither strategy is appropriate in a dense mapping context as drift in odometry estimation before loop detection may cause unrelated areas of the map to completely overlap in space. This issue also arises in determining connectivity of the graph. Applying sampling and connectivity strategies that are only spatially aware can result in links between completely unrelated points in the map, as shown in Figure 11. The effects of applying a nearest neighbour strategy are visualised in Figure 12. For this reason we derive a sampling and connectivity strategy that exploits the camera pose graph for deformation graph construction and connection. The method is computationally efficient and incremental, enabling real-time execution. Our sampling strategy is listed in Algorithm 4 in Appendix B.

We connect deformation graph nodes returned by our sampling strategy in a sequential manner, following the temporal order of the pose graph itself. That is to say our set of graph nodes  $N$  is ordered. We sequentially connect nodes up to a value  $k$ . We use  $k = 4$  in all of our experiments. For example, a node  $l$  will be connected to nodes  $(l \pm 1, l \pm 2)$ . We show  $k = 2$  connectivity in Figure 11. Note the connectivity of end nodes which maintains  $k$ -connectivity.



**Figure 12:** From left to right; (i) Highly distorted map produced when a naïve nearest neighbour sampling and connectivity strategy is used; In this example, parts of the floor close the point of loop closure have been associated with the nearby window through the deformation graph. When optimised, these parts of the scene attempt to “stick together”, drastically distorting the surrounding geometry. (ii) Non-distorted map loop closure using our proposed sampling and connectivity strategy. When the deformation graph is intelligently constructed across the map using our scheme, incorrect surface association problems as shown on the left are avoided.

### 4.3.3 Incremental Vertex Weighting

Each vertex  $\mathbf{v}$  has a set of influencing nodes in the deformation graph  $\mathcal{N}(\mathbf{v})$ . The deformed position of a vertex is given by Sumner et al. (2007):

$$\hat{\mathbf{v}} = \sum_{k \in \mathcal{N}(\mathbf{v})} w_k(\mathbf{v}) \left[ N_k^{\mathbf{R}}(\mathbf{v} - N_k^{\mathbf{g}}) + N_k^{\mathbf{g}} + N_k^{\mathbf{t}} \right] \quad (30)$$

where  $w_k(\mathbf{v})$  is defined as (all  $k$  summing to 1):

$$w_k(\mathbf{v}) = (1 - \|\mathbf{v} - N_k^{\mathbf{g}}\|_2 / d_{max})^2 \quad (31)$$

Here  $d_{max}$  is the Euclidean distance to the  $k+1$ -nearest node of  $\mathbf{v}$ . In previous work based on this technique the sets  $\mathcal{N}(\mathbf{v})$  for each vertex are computed in batch using a  $k$ -nearest neighbour technique. Again, being based on spatial constraints alone this method fails in the example shown in Figure 11. To overcome this issue we derive an algorithm that assigns nearest neighbour nodes to each vertex using a greedy back-traversal of the sampled pose graph nodes.

Referring back to Figure 10 and Section 2.3.2, we recall that each pose that causes a volume shift has an associated set of vertices contained within a cloud slice. We can exploit the inverse mapping of this association to map each vertex onto a single pose in the pose graph. However, the associated pose is at least a distance of  $\frac{v_d}{2}$  away from the vertex, which is not ideal for the deformation. In order to pick sampled pose graph nodes for each vertex that are spatially and temporally optimal, we use the closest sampled pose to the associated cloud slice pose as a starting point to traverse back through the sampled pose graph nodes to populate a set of candidate nodes. From these candidates the  $k$ -nearest neighbours of the vertex are chosen. We list the algorithm for this procedure in Algorithm 5 in Appendix B and provide a visual example in Figure 10.

The per-vertex node weights can be computed within the back-traversal algorithm, which itself can be carried out incrementally online while the frontend volume shifting component provides new cloud slices. The ability to avoid computationally expensive batch steps for deformation graph construction and per-vertex weighting by using incremental methods is the key to allowing low latency online map optimisation at any time.

## 4.4 Optimisation

On acceptance of a loop closure constraint as described in Section 4.2 we perform two optimisation steps, firstly on the pose graph and secondly on the dense vertex map. The pose graph optimisation provides the measurement constraints for the dense map deformation optimisation in place of user specified constraints that were necessary in the original embedded deformation approach. Pose graph optimisation is carried out using the iSAM framework (Kaess et al. (2008)). We benefit from the incremental sparse linear algebra representation used internally in iSAM, such that execution time is reasonable in terms of online operation.

### 4.4.1 Map Deformation

Sumner et al. (2007) define three cost functions over the deformation graph and user constraints to optimise the set of affine transformations over all graph nodes  $N$ . The first maximises rigidity in the deformation:

$$E_{rot} = \sum_l \left\| N_l^{\mathbf{R}^T} N_l^{\mathbf{R}} - \mathbf{I} \right\|_F^2 \quad (32)$$

Where Equation 32 is the alternative Frobenius-norm form provided by Chen et al. (2012). The second is a regularisation

term that ensures a smooth deformation across the graph:

$$E_{reg} = \sum_I \sum_{n \in \mathcal{N}(N_I)} \|N_I^R(N_n^g - N_I^g) + N_I^g + N_I^t - (N_n^g + N_n^t)\|_2^2 \quad (33)$$

The third is a constraint term that minimises the error on a set of user specified vertex position constraints  $Q$ , where a given constraint  $Q_p \in \mathbb{R}^3$  and  $\phi(\mathbf{v})$  is the result of applying Equation 30 to  $\mathbf{v}$ :

$$E_{con} = \sum_p \|\phi(\mathbf{v}) - Q_p\|_2^2 \quad (34)$$

We link the optimised pose graph to the map deformation through the  $E_{con}$  cost function. With  $P$  being the pose graph (composed of rotations and translations  $\mathbf{R}_i$  and  $\mathbf{t}_i$ ) before loop constraint integration we set  $P'$  to be the optimised pose graph returned from iSAM. We then add each of the camera pose translations to the deformation cost as if they were user specified vertex constraints, redefining Equation 34 as:

$$E_{con_p} = \sum_i \|\phi(\mathbf{t}_i) - \mathbf{t}'_i\|_2^2 \quad (35)$$

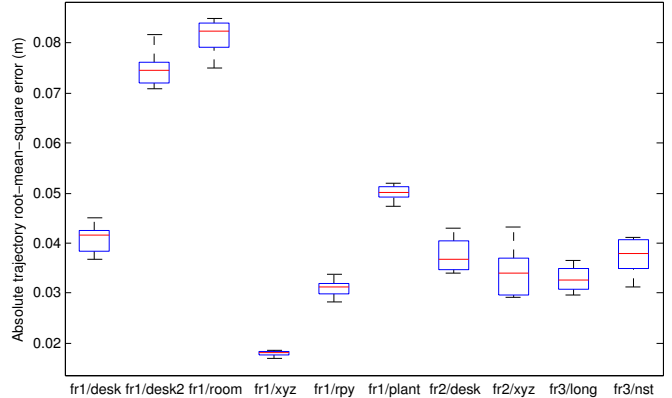
A uniform constraint distribution across the surface obtained from this parameterisation aids in constraining both surface translation and orientation. However at some points the surface orientation may not be well constrained. In order to overcome this issue we add additional vertex constraints between the unoptimised and optimised 3D back-projections of each of the matched inlier SURF keypoints detected in Section 4.2, where  $P_i(\mathbf{R}_i$  and  $\mathbf{t}_i)$  is the camera pose of the matched loop closure frame:

$$E_{surf} = \sum_q \|\phi((\mathbf{R}_i G_q) + \mathbf{t}_i) - ((\mathbf{R}'_i G_q) + \mathbf{t}'_i)\|_2^2 \quad (36)$$

The final total cost function is defined as:

$$w_{rot} E_{rot} + w_{reg} E_{reg} + w_{con_p} E_{con_p} + w_{surf} E_{surf} \quad (37)$$

With  $w_{rot} = 1$ ,  $w_{reg} = 10$ ,  $w_{con_p} = 100$  and  $w_{surf} = 100$ , we minimise this unnormalised cost function using the iterative Gauss-Newton algorithm choosing weighting values in line with those used in Sumner et al. (2007). The optimisation consistently converges to a satisfactory result with these weights, similar to the findings of Chen et al. (2012). As highlighted in previous work, the Jacobian matrix in this problem is sparse, enabling the use of sparse linear algebra libraries for efficient optimisation. We use the CHOLMOD library to perform sparse Cholesky factorisation and efficiently solve the system (Davis and Hager (1999)). We then apply the optimised deformation graph  $N$  to all vertices over all cloud slices  $C$  in parallel across multiple CPU threads. As discussed in Section 2.3.2 we compute an incremental mesh surface representation of the cloud slices as they are produced by the frontend. The incremental mesh can be deformed by applying the deformation graph to its vertices. In our experience an incremental mesh typically contains more minuscule holes than



**Figure 13:** Boxplot of the ATE RMSE in metres per sequence evaluated. In each box the red central line is the median, the box edges the 25th and 75th percentiles and the whiskers extend to the minimum and maximum estimates. Each dataset was ran ten times to account for the randomness induced by the place recognition system in Section 4.2.

a batch mesh, which in path planning is functionally almost identical but less visually appealing. In all results we show the batch mesh computed over the set of optimised vertices.

## 5 Evaluation

We evaluate our system both quantitatively and qualitatively in terms of trajectory estimation, surface reconstruction and computational performance. We processed a combined total of over 79,000 unique RGB-D frames in our evaluation.

### 5.1 Trajectory Estimation

To evaluate the accuracy of our camera trajectory estimation we present results on the widely used RGB-D benchmark of Sturm et al. (2012). This benchmark provides synchronised ground truth poses for an RGB-D sensor moved through an environment, captured with a highly precise motion capture system. We evaluated multiple runs over ten datasets with quantitative results shown in Table 1 and a boxplot shown in Figure 13. We use the absolute trajectory (ATE) root-mean-square error metric (RMSE) to evaluate our system, which measures the root-mean-square of the Euclidean distances between all estimated camera poses and the ground truth poses associated by timestamp (Sturm et al. (2012)).

Consistent performance is achieved on all sequences evaluated, with a notably higher error on the fr1/desk2 and fr1/room datasets. This can be explained by the high average angular velocity on these sequences which causes motion blur, increases the effect of rolling shutter and violates the assumption of projective data association. From the results it can be seen that a higher RMSE is correlated with a high average angular velocity. Provided there is a low standard deviation in frame rate and good overlap between successive frames a strong trajectory estimate is achievable. Figure 14 shows two dimensional plots of the differences between the

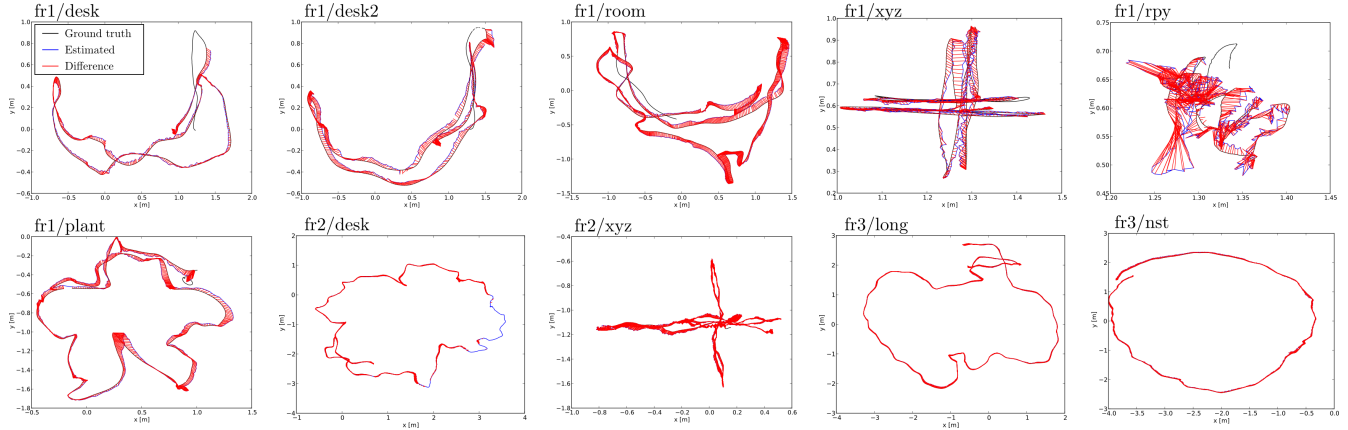


Figure 14: Two dimensional plot of estimated trajectories versus ground truth trajectories on evaluated sequences.

Dataset	RMSE (m)	Median (m)	Max (m)	$\bar{\omega}$ ( $^{\circ}$ -1)
fr1/desk	0.0407	0.0352	0.0905	23.33
fr1/desk2	0.0747	0.0639	0.2309	29.31
fr1/room	0.0813	0.0739	0.2511	29.88
fr1/xyz	0.0180	0.0155	0.0392	8.92
fr1/rpy	0.0311	0.0213	0.0991	50.15
fr1/plant	0.0500	0.0425	0.1148	27.89
fr2/desk	0.0376	0.0315	0.0879	6.34
fr2/xyz	0.0341	0.0234	0.0979	1.72
fr3/long	0.0329	0.0297	0.0698	10.19
fr3/nst	0.0372	0.0335	0.0735	7.43

Table 1: Statistics on ATE on evaluated datasets. Trajectory values are in metres as the mean over ten runs of each dataset. The mean angular velocity is given as  $\bar{\omega}$  in degrees per second, retrieved from the dataset specifications.

Dataset	Ours (m)	DVO (m)	RGB-D (m)	MRS (m)
fr1/desk	0.037	0.021	0.023	0.043
fr1/desk2	0.071	0.046	0.043	0.049
fr1/room	0.075	0.053	0.084	0.069
fr1/xyz	0.017	0.011	0.014	0.013
fr1/rpy	0.028	0.020	0.026	0.027
fr1/plant	0.047	0.028	0.091	0.026
fr2/desk	0.034	0.017	0.057	0.052
fr2/xyz	0.029	0.018	0.008	0.020
fr3/long	0.030	0.035	0.032	0.042
fr3/nst	0.031	0.018	0.017	-

Table 2: Comparison of ATE RMSE on evaluated datasets and SLAM systems. All units given are in metres. MRS was unable to produce an estimate on the fr3/nst dataset.

estimated trajectories and the ground truth trajectories. In all real world datasets evaluated in this paper the auto exposure and auto white balance features of the RGB-D camera were enabled.

### 5.1.1 Comparative Evaluation

We compare the trajectory estimation performance of our system to three recent state-of-the-art visual SLAM systems, DVO SLAM of Kerl et al. (2013), RGB-D SLAM of Endres et al. (2012) and multi-resolution surfel maps (MRS) of Stückler and Behnke (2013). Table 2 summarises the results, where our values represent the best estimate over ten runs. From these we can see the performance of our system is comparable to other leading approaches, where performance of each algorithm is typically within no more than 3cm in total RMSE. We acknowledge the strong performance of the DVO SLAM system of Kerl et al. (2013) in trajectory estimation and perform a further comparison with their system in terms of reconstruction accuracy and larger trajectories in Section 5.2.2. We also provide a small comparison of results between our system and benchmark results provided by Meilland and Comport (2013) from their unified keyframe

Dataset	Ours (m)	Unified (m)
fr1/desk	0.031	0.018
fr2/desk	0.028	0.093
fr1/room	0.068	0.144
fr2/large_no_loop	0.256	0.187

Table 3: Comparison of ATE Median error on evaluated datasets and SLAM systems. All units given are in metres.

SLAM system in Tables 3 and 4, again showing comparable performance (using their chosen metric of ATE Median and Max error, as opposed to RMSE). Note that the results on the fr2/large\_no\_loop dataset are taken from our previous work, Whelan et al. (2013a).

Dataset	Ours (m)	Unified (m)
fr1/desk	0.078	0.066
fr2/desk	0.079	0.116
fr1/room	0.231	0.339
fr2/large_no_loop	0.878	0.317

Table 4: Comparison of ATE Max error on evaluated datasets and SLAM systems. All units given are in metres.



## 5.2 Surface Reconstruction

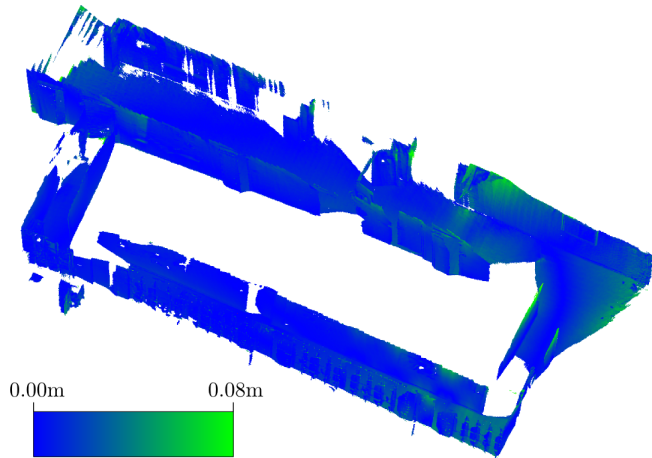
We present a number of quantitative and qualitative results on evaluating the surface reconstructions produced by our system. In our experience a high score on a camera trajectory benchmark does not always imply a high quality surface reconstruction due to the frame-to-model tracking component of the system. In previous work we found that although other methods for camera pose estimation may score better on benchmarks, the resulting reconstructions are not as accurate if frame-to-model tracking is not being utilised (Whelan et al. (2013a)). We evaluate seven different datasets captured in a handheld fashion across a wide range of environments, demonstrating the viability of our system for use over large scale trajectories both indoors and outdoors (within sensing limitations) and across multiple floors. It should be noted that it is technically possible for self-intersection to occur in the surface upon deformation. We have found this to be quite rare in practice as most deformations are quite smooth and do not deform the map in erratic ways (visualised in Multimedia Extension 1). This aspect of the algorithm is one of the trade offs made in favor of computational performance.

### 5.2.1 Comparison to 2-pass Optimisation

In order to evaluate the accuracy of the deformation process we compare the resulting maps produced when a 2-pass approach is taken versus a single pass approach with a deformation for map correction. The 2-pass approach involves the following steps;

1. Build a pose graph with a camera pose for every frame.
2. Detect visual loop closures using the method described in Section 4.2.
3. At the end of the dataset, optimise the camera pose graph taking loop closure constraints into account.
4. Rerun the dataset using the optimised pose graph in place of the visual odometry frontend.

From here we can compare the two maps to determine a measure of similarity. This presents an interesting question as although the pose graphs for both the 2-pass and deformation-based maps are identical, the maps themselves may differ slightly due to the fact the 2-pass approach gives up frame-to-model registration on the second pass where the frustum-volume intersection may also slightly change. This means there will not be any reliable 1-to-1 point correspondences between the maps. For this reason we measure the map similarity by the residual error of a dense ICP-based registration of the maps. Given that both maps lie in the global coordinate frame we can iteratively minimise nearest neighbour point-wise correspondences between the two maps using standard point-to-plane ICP. This allows us to account for a small rigid transformation error between the two maps. We measure the remaining root-mean-square residual error between point correspondences as the residual similarity error between the two



**Figure 15:** Heatmap showing the difference between the deformed reconstruction and 2-pass reconstruction of the Indoor dataset. Blue indicates no error and scales to pure green indicating an deviation of 0.08m.

maps. Table 5 lists statistics on the seven evaluated datasets including the 2-pass residual registration error as well as the same error computed on maps deformed with a subsampled pose graph, which we discuss in Section 5.3. It is clear that the deformation approach brings the map into strong alignment with the 2-pass output, with only a few millimetres in difference. This can be seen in Figure 15. Multimedia Extension 1 shows the map correcting deformation occurring for the Indoors and Two floors datasets, as well as flythroughs of the final meshes. Images of all datasets are provided in Figures 22-28 in the Appendix. The Apartment dataset has a notably higher error than the other sequences, owing to the complexity of the trajectory and scene. However observing the reconstruction in Figure 28 it can be seen that a high quality map is still achieved.

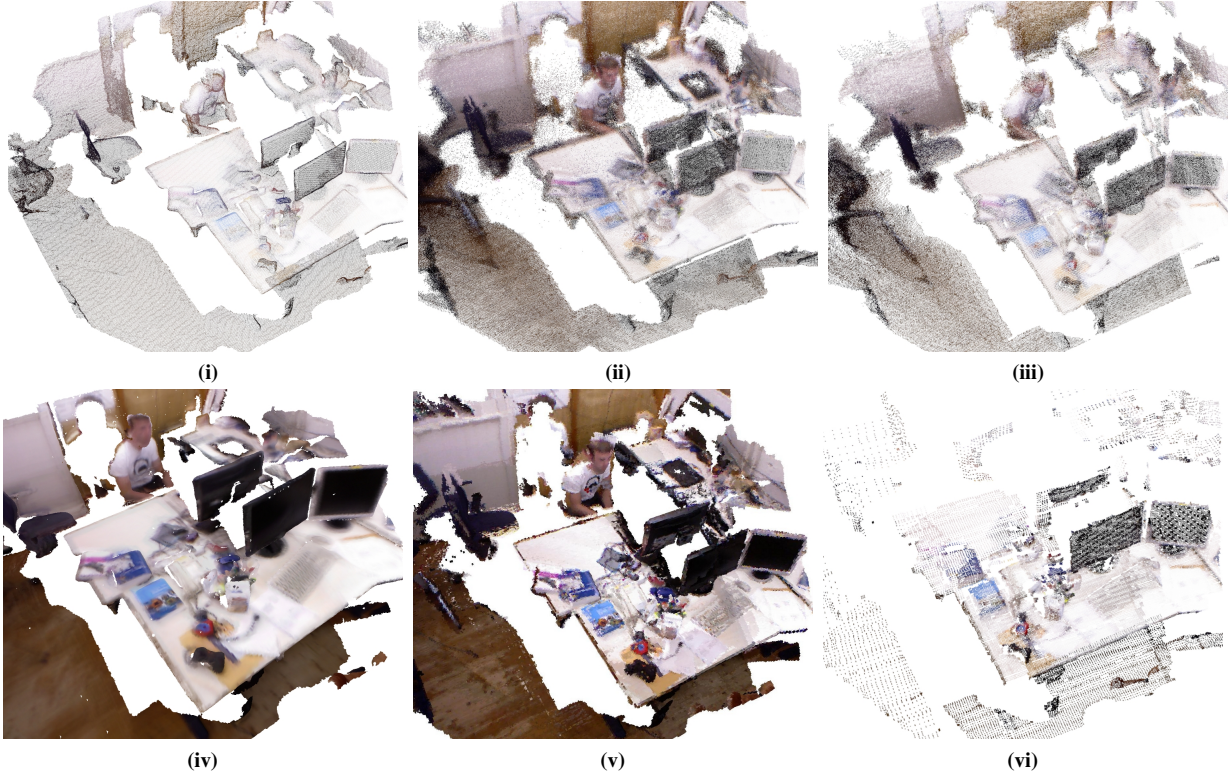
### 5.2.2 Surface Reconstruction Comparison

In Figure 16 we present a comparison of the reconstructions produced by each of the systems evaluated in Section 5.1.1 on the fr1/xyz data. From this qualitative comparison it is evident that our approach benefits greatly from the use of a fused volumetric frontend, removing substantial noise from the reconstruction and producing a much cleaner model than other approaches. While the output from RGB-D SLAM, MRSSMap and DVO SLAM can be fed through a signed distance fusion pipeline to produce a similar output, this would strictly be a post-processing step that is not required by our own system to produce such reconstructions.

We also compare our reconstruction results on all of our seven evaluated datasets to the output produced by the open source DVO SLAM system of Kerl et al. (2013), using the provided default configuration parameters. DVO SLAM performance statistics are listed in Table 6. In all datasets the DVO SLAM system frontend executed at 30Hz. The final

Dataset	Length(m)	$v_d$ (m)	$d_p$ (m)	Vertices	Volume(m <sup>3</sup> )	2-pass(mm)	2-pass fast(mm)	Figure
Coffee	30.18	4.5	0.4	909,422	7,993	1.2	1.8	22
Indoors	49.57	7	0.4	1,603,116	21,918	2.7	4.9	23
Garden	71.49	6	0.8	2,418,331	28,340	2.1	2.5	24
Outdoors	152.05	6	0.8	2,961,966	34,711	2.3	8.8	25
Two floors	173.88	6	0.8	4,016,273	47,066	2.1	8.0	26
In/outdoors	317.95	6	0.8	5,985,669	70,145	2.8	7.5	27
Apartment	61.27	2.5	0.3	6,205,222	30,299	19.0	20.4	28

**Table 5:** Statistics on seven handheld datasets captured over a wide variety of environments using our approach.



**Figure 16:** Comparison of reconstructions on the fr1/xyz dataset; (i) Point cloud reconstruction with our approach showing a smooth surface reconstruction. (ii) Reprojected keyframe reconstruction from RGB-D SLAM, showing a noisy surface with quantization effects. (iii) Reprojected keyframe reconstruction from DVO SLAM again showing a noisy surface with quantization effects. (iv) Triangular mesh reconstruction with our approach. (v) OctoMap (Hornung et al. (2013)) reconstruction from RGB-D SLAM, while in this form useful for motion planning, appears very jagged and is quantized to the nearest voxel. (vi) Point cloud sampled at highest resolution from surfel map with MRSMap showing an evident discretization effect.

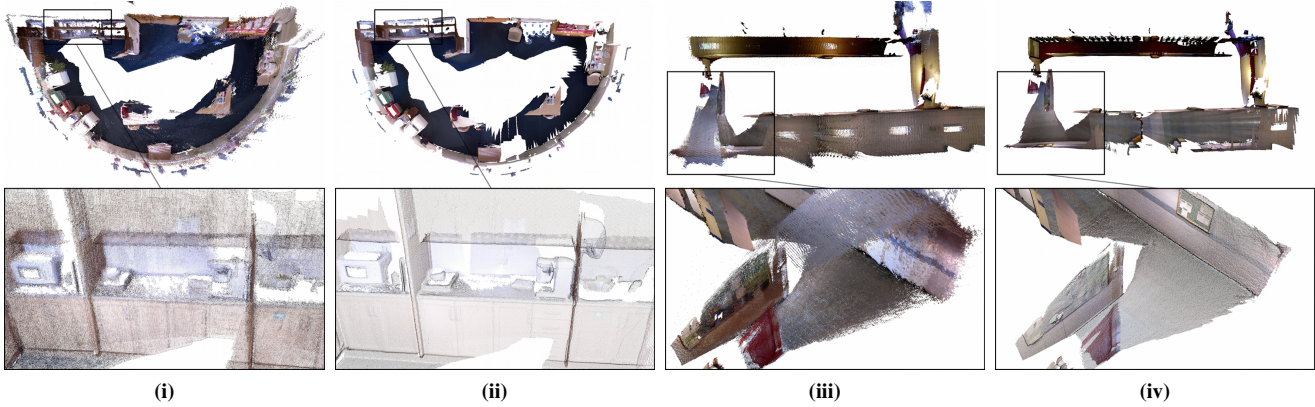
pose graph optimisation and additional keyframe loop closure search time is listed in the “Post-processing” column. The DVO SLAM system uses no specific method for map reconstruction and must rely on point cloud reprojection of raw RGB-D keyframes to reconstruct the map after the pose graph has been optimised. This results in many redundant and repeated points in the map. To remedy this problem we apply a voxel grid downsampling filter (as mentioned in Section 2.3.2) with a resolution of 1cm to the output keyframe vertices to keep the map size tractable. These numbers are listed in the “Vertices” and “Vertices (filtered)” columns. As listed the system successfully reconstructs the Coffee and Indoors datasets, however in contrast to our approach post-processing time of between 7 and 31 seconds is required to optimise the final pose graph and resolve any additional keyframe loop

asures (where our system does not require any post-processing or final batch steps). Failure to detect loop closures results in failed reconstructions on the Garden, Outdoors, Two floors, In/outdoors and Apartment datasets which could perhaps be remedied by using a bag-of-words visual features-based approach similar to ours (Galvez-Lopez and Tardos (2011)) or indeed as suggested by Kerl et al. (2013) the FAB-MAP algorithm (Cummins and Newman (2010)). Camera pose estimation failures were also encountered in the Outdoors, Two floors, In/outdoors and Apartment datasets, particularly in regions of the sequences which were mostly planar or had strong visual aliasing, such as staircases. From these results and those listed in Table 2 we observe that the method for detecting loop closures used by Kerl et al. (2013) is very strong in small sized environments but scales poorly as the explored



Dataset	Vertices	Vertices (filtered)	Post-processing (s)	Verdict	Figure
Coffee	34,813,313	3,925,307	7.41	Successful	17 (i)
Indoors	56,927,008	15,960,983	31.65	Successful	17 (iii)
Garden	145,054,728	15,385,263	159.52	Loop Closure Failure	N/A
Outdoors	104,948,878	9,106,848	190.76	Loop Closure & Tracking Failure	N/A
Two floors	237,308,674	26,724,533	841.93	Loop Closure & Tracking Failure	N/A
In/outdoors	275,096,207	18,836,984	5501.58	Loop Closure & Tracking Failure	N/A
Apartment	83,102,006	6,068,711	74.79	Loop Closure & Tracking Failure	N/A

**Table 6:** Statistics on seven handheld datasets captured over a wide variety of environments processed using the DVO SLAM system.



**Figure 17:** From left to right; (i) DVO SLAM keyframe reprojection of the Coffee dataset. The surfaces are notably noisy and quantization effects are evident. (ii) Our reconstruction of the Coffee dataset, showing a smooth uniform reconstruction. (iii) DVO SLAM keyframe reprojection of the Indoors dataset. Again surfaces are very noisy and highly quantized. In contrast to our reconstruction, the ceiling has been mapped in most of the sequence, however being quite distant from the sensor suffers badly from discretization effects. (iv) Our reconstruction of the Indoors dataset. The ceiling has not been reconstructed in this sequence since the configuration of the TSDF volume size caused it to fall outside of the area of reconstruction. This is however easily remedied by modifying the relative parameterisation of the volume with respect to the sensor, similar to the dynamic cube positioning technique we discussed in Section 2.4.

area size grows, both in terms of accuracy and computational performance (embodied in the consistently increasing post-processing time).

In Figure 17 we qualitatively compare the reconstruction quality of our approach versus the maps produced by DVO SLAM on the Coffee and Indoors datasets. For clarity we compare vertices only as DVO SLAM provides no method for mesh surface reconstruction. These results show that the reconstruction produced by our approach is much smoother and contains significantly fewer redundant vertices. Additionally, there are no raw RGB-D point cloud quantization effects in our reconstructions. The reprojection principle taken to producing the map from DVO SLAM keyframes does result in entire frame back-projection which produces a “fuller” looking map, however far away points in current generation RGB-D sensors are known to be extremely noisy and highly inaccurate (Khoshelham and Elberink (2012)).

### 5.2.3 Surface Ground Truth

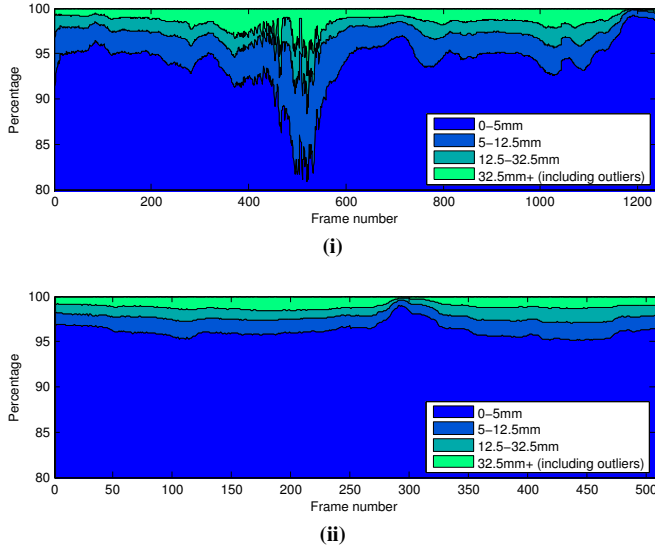
We evaluate the surface reconstruction quality of our approach quantitatively using synthetic data produced in an identical manner to the datasets created by Handa et al. (2012). Each dataset contains 30Hz RGB-D frames from a camera placed in a synthetic office environment. The camera



**Figure 18:** Mesh reconstruction of the first synthetic dataset. Note that the rough triangulation of parts of the chairs is due to a poor viewing angle throughout the sequence.

trajectories were generated from real world data which was previously ran through our visual odometry frontend. Given that the datasets were produced using a procedural raytracing process (using POVRay), there is no actual surface to compare against. However, each RGB-D frame does have ground truth depth information which we compare against. For each frame in a dataset we compute a histogram of the per depth pixel  $L_1$ -norm error between the ground truth depth map and





**Figure 19:** Temporal histograms of predicted depth versus ground truth depth on synthetic datasets. A frame from the dip in accuracy around the center of the first dataset is shown in Figure 20 (i) while a frame from the peak in accuracy in the center of the second dataset is shown in Figure 20 (ii).

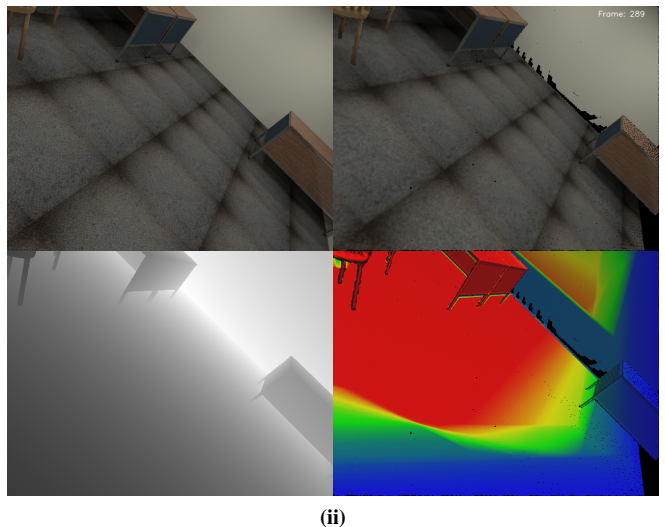
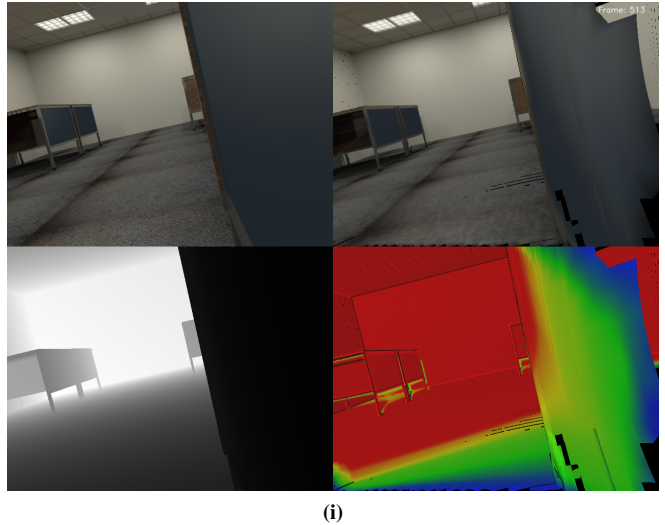
the predicted surface depth map raycast from the TSDF, normalising by the number of valid pixels before aligning all histograms into a two dimensional area plot. We evaluated two synthetic datasets of the same scene with different camera motions. The temporal error histograms are shown in Figure 19 while frames from each dataset are shown in Figure 20. Overall the synthetic surfaces are reconstructed very well, however occasional raycasting artifacts (particularly around the edges of objects and on nearby surfaces) can hinder the reconstruction quality score, as in the first dataset. These artifacts occur due to the use of a fixed step size during ray casting and the absence of any special method to render smooth edges, both for performance reasons. Observing the final reconstruction in Figure 18 it is clear that the slight dip in accuracy did not effect the reconstruction quality by any significant amount. Typically around 95% of the estimated depth of the surface is within 5mm of ground truth.

### 5.3 Computational Performance

We evaluate the computational performance of both the frontend and backend of the system. The evaluation platform was a standard desktop PC running Ubuntu 12.04 with an Intel Core i7-3960X CPU at 3.30GHz, 16GB of RAM and an nVidia GeForce 680GTX GPU with 2GB of memory.

#### 5.3.1 Frontend Performance

To evaluate the performance of the frontend (including volume integration, camera pose estimation, volume raycasting and volume shifting, essentially all teal colored function blocks in Figure 9) we provide frame processing timing re-



**Figure 20:** One frame from each surface ground truth evaluation dataset. Each shows in clockwise order the ground truth RGB, predicted RGB, predicted surface phong shaded colored by voxel weight and ground truth depth map. From top to bottom; (i) Here raycasting artifacts are visible in the predicted surface in the bottom right causing a high error in the evaluation; This is evident particularly in the top right-hand corner of the frame where the wall is visible through the side of the desk. (ii) Overall the surface is being well estimated and there are no raycasting artifacts.

$m_s$	Avg (ms)	Min (ms)	Max (ms)	StdDev (ms)
1	34.15	25.93	41.58	3.30
2	32.21	25.63	39.29	3.14
4	31.08	25.38	39.02	2.77
8	30.57	25.42	37.44	2.48
<b>16</b>	<b>29.94</b>	<b>24.97</b>	<b>37.25</b>	<b>2.26</b>
32	30.26	25.33	40.30	2.39
64	30.49	25.06	43.95	2.73

**Table 7:** Computational performance of the volumetric fusion thread on the fr1/desk dataset. The shifting threshold  $m_s$  is given in voxels while the frame processing timings are given in milliseconds. Highlighted is the optimal choice based on execution time.

sults on the fr1/desk sequence comparing different choices of the  $m_s$  parameter discussed in Section 2.3. This parameter affects the frequency and size of each volume shift, which in turn affects frontend performance. Results are shown in Table 7. A shifting threshold of 16 voxels was found to be optimal, providing the best computational performance with an average frame rate comfortably below the frame rate of the sensor (30Hz) and with minimal spikes in execution time.

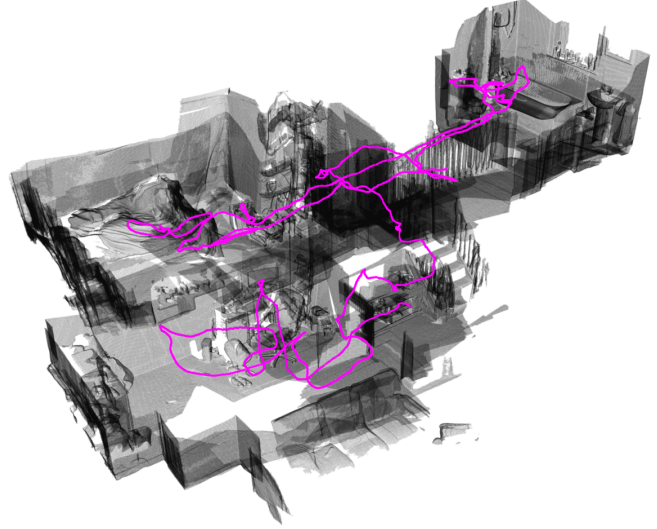
### 5.3.2 Backend Performance

We quantify the computational performance of the backend in the context of an online real-time SLAM system by measuring the latency of the system. That is, how long it takes for 1) a loop closure to be recognised when one is encountered and 2) map correction to be completed. Table 8 shows execution time and latency statistics on our test platform for the first six datasets, while Table 10 shows performance statistics on the Apartment dataset. We also experimented with subsampling the pose graph used in the iSAM-based pose graph optimisation by the same sampling metric used in Algorithm 4. This affects the number of poses used in the final pose graph optimisation and the number of points available to constrain the map deformation in Equation 35. Our results (shown in Tables 9 and 11) show that using a subsampled pose graph (akin to using only keyframes) instead of an every frame pose graph reduces execution time (and therefore latency) by up to almost an order of magnitude in some cases, while only mildly affecting map quality (quantified as “2-pass fast” in Table 5). As expected the appearance-based frontend scales very well over hundreds of metres while the backend is capable of correcting millions of vertices for global consistency in only 1-3 seconds. The results presented in Tables 10 and 11 demonstrate the capability of our approach to deal with complex trajectories with multiple loop closures. This is further highlighted by the plot of the camera trajectory on the Apartment dataset shown in Figure 21.

Multimedia Extension 2 shows the entirety of the In/outdoors dataset running in real-time including the two online loop closures while Multimedia Extension 3 shows the Apartment dataset. Note that in these videos the vertex count is higher due to the weight-based filtering mentioned in Section 2.3.2 being disabled, resulting in more extracted vertices from the TSDF slices.

## 6 Conclusion

In this paper we have presented a real-time dense SLAM system which makes use of a dense every-frame volumetric fusion frontend for camera pose estimation and surface reconstruction in combination with a non-rigid map deformation backend to correct the mapped dense surface upon loop closure. We have provided an extensive evaluation, both quantitatively and qualitatively on common benchmarks and our own datasets demonstrating the system’s ability to produce large scale dense globally consistent maps in real-time.



**Figure 21:** Camera trajectory plot within the Apartment dataset, showing the “loopy” path the camera took through the environment.

One limitation in our system is the reliance on projective data association for camera pose estimation which limits the kinds of motion that our visual odometry frontend can handle. However this restriction works in our favour as with increased camera motion comes increased motion blur and rolling shutter effects. Approaches exist to correct for such effects in real-time such as that of Meilland et al. (2013), however this would cause an increased computational requirement when aggregated to any existing system.

Our current implementation does not support the reintegration of areas of the map which are revisited into the volumetric fusion frontend. This results in aliasing in areas that receive multiple passes. However representing the surface as a set of cloud slices maintains spatiotemporal information about the map which can be used for change detection, scene differencing or even the merging of cloud slices from multiple passes. Reintegration or re-fusing of the mesh-based map in real-time is a challenging problem due to the sheer volume of data. Some existing approaches discussed in Section 1.1 support this but lack a means for correcting for drift or global consistency online in real-time. Commonly adopted space efficient data structures need to be fully restructured upon large updates to the map, which in turn would hinder real-time performance greatly. Other discussed approaches are either offline, or sacrifice local surface connectivity to achieve surface refusing. Real-time large scale dense fused 3D reconstruction which supports online drift correction, provides a globally consistent representation of the map at any time and allows map re-use and re-fusing is a challenging problem which we aim to address in our future work. We also plan to research new methods for estimating camera pose uncertainty and scalability over hundreds of metres.

Quantities	Datasets						
	Coffee	Indoors	Garden	Outdoors	Two floors	In/outdoors(1)	In/outdoors(2)
DBoW images	280	301	658	1171	1584	1662	2706
Poses	1544	2993	8634	5240	12952	17306	25586
Nodes	58	55	72	178	191	211	364
Vertices	932,056	1,352,919	2,256,475	2,805,083	3,896,281	3,560,994	5,867,125
Process	Timings (ms)						
Frontend	465	602	622	587	657	521	543
iSAM	257	510	1412	1299	3326	4386	6545
Deformation	266	390	1112	928	2197	3040	4473
Total latency	988	1502	3146	2814	6180	7947	11561

**Table 8:** Computational performance statistics on six datasets using an every frame pose graph. Quantities shown are at the moment of loop closure. The In/outdoors dataset contains two looping points which are both listed.

Quantities	Datasets						
	Coffee	Indoors	Garden	Outdoors	Two floors	In/outdoors(1)	In/outdoors(2)
DBoW images	277	305	672	1173	1593	1713	2782
Poses	283	307	674	1186	1594	1716	2783
Nodes	52	49	68	167	181	196	339
Vertices	943,721	1,371,560	2,246,028	2,841,135	3,904,113	3,569,842	5,850,152
Process	Timings (ms)						
Frontend	488	589	651	597	467	540	793
iSAM	46	67	110	288	378	271	1140
Deformation	110	105	170	377	381	148	842
Total latency	644	761	931	1262	1226	959	2775

**Table 9:** Computational performance statistics on six datasets using a subsampled pose graph. Quantities shown are at the moment of loop closure. The In/outdoors dataset contains two looping points which are both listed.

Loop number	Apartment dataset with full pose graph				
	1	2	3	4	5
DBoW images	119	526	708	982	1428
Poses	367	1638	2163	2824	3937
Nodes	14	61	80	105	165
Vertices	492,960	2,792,446	3,800,812	4,482,186	6,296,542
Process	Timings (ms)				
Frontend	807	858	1596	703	604
iSAM	29	202	277	230	648
Deformation	51	336	425	425	932
Total latency	887	1396	2298	1358	2184

**Table 10:** Computational performance statistics on the Apartment dataset using an every frame pose graph. Quantities shown are at the moment of loop closure.

Loop number	Apartment dataset with subsampled pose graph				
	1	2	3	4	5
DBoW images	119	529	708	982	1430
Poses	123	531	715	988	1433
Nodes	13	59	77	100	157
Vertices	492,718	2,791,445	3,799,464	4,490,170	6,295,379
Process	Timings (ms)				
Frontend	789	868	1557	789	593
iSAM	19	64	93	88	235
Deformation	31	181	252	285	508
Total latency	839	1113	1902	1162	1336

**Table 11:** Computational performance statistics on the Apartment dataset using a subsampled pose graph. Quantities shown are at the moment of loop closure.

## 7 Acknowledgements

Research presented in this paper was funded by a Strategic Research Cluster grant (07/SRC/I1168) by Science Foundation Ireland under the Irish National Development Plan, the Embark Initiative of the Irish Research Council, ONR grants N00014-10-1-0936, N00014-11-1-0688, N00014-12-1-0093, N00014-12-10020 and NSF grant IIS-1318392. The authors would like to thank Guillaume Gales, Richard H. Middleton and Ross Finman for their input and discussion and also Ankur Handa for his synthetic surface ground truth datasets.

## References

- Audras, C., Comport, A. I., Meilland, M., and Rives, P. (2011). Real-time dense RGB-D localisation and mapping. In *Australian Conf. on Robotics and Automation*, Monash University, Australia.
- Bylow, E., Sturm, J., Kerl, C., Kahl, F., and Cremers, D. (2013). Real-time camera tracking and 3d reconstruction using signed distance functions. In *Robotics: Science and Systems Conference (RSS)*.
- Canelhas, D. R., Stoyanov, T., and Lilienthal, A. J. (2013). SDF Tracker: A parallel algorithm for on-line pose estimation and scene reconstruction from depth images. In *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems, IROS*, Tokyo, Japan.
- Chen, J., Bautembach, D., and Izadi, S. (2013). Scalable real-time volumetric surface reconstruction. In *SIGGRAPH 2013*, Anaheim, CA, USA. ACM.
- Chen, J., Izadi, S., and Fitzgibbon, A. (2012). KinÊtre: animating the world with the human body. In *Proceedings of the 25th annual ACM symposium on User interface software and technology, UIST '12*, pages 435–444, New York, NY, USA. ACM.
- Comport, A., Malis, E., and Rives, P. (2007). Accurate Quadri-focal Tracking for Robust 3D Visual Odometry. In *IEEE International Conference on Robotics and Automation, ICRA'07*, Rome, Italy.
- Cummins, M. and Newman, P. (2010). Invited Applications Paper FAB-MAP: Appearance-Based Place Recognition and Mapping using a Learned Visual Vocabulary Model. In *27th Intl Conf. on Machine Learning (ICML2010)*.
- Davis, T. A. and Hager, W. W. (1999). Modifying a sparse Cholesky factorization. *SIAM J. Matrix Anal. Appl.*, 20(3):606–627.
- Deutsch, P. and Gailly, J.-L. (1996). Zlib compressed data format specification version 3.3.
- Endres, F., Hess, J., Engelhard, N., Sturm, J., Cremers, D., and Burgard, W. (2012). An evaluation of the RGB-D SLAM system. In *Proc. of the IEEE Int. Conf. on Robotics and Automation (ICRA)*, St. Paul, MA, USA.
- Fischler, M. A. and Bolles, R. C. (1981). Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395.
- Galvez-Lopez, D. and Tardos, J. D. (2011). Real-time loop detection with bags of binary words. In *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ Int. Conf. on*, pages 51–58.
- Handa, A., Newcombe, R., Angeli, A., and Davison, A. (2012). Real-time camera tracking: When is high frame-rate best? In *ECCV 2012*, volume 7578 of *Lecture Notes in Computer Science*, pages 222–235.
- Henry, P., Fox, D., Bhowmik, A., and Mongia, R. (2013a). Patch volumes: Multiple fusion volumes for consistent RGB-D modeling. In *RSS Workshop on RGB-D: Advanced Reasoning with Depth Cameras*, Berlin, Germany.
- Henry, P., Fox, D., Bhowmik, A., and Mongia, R. (2013b). Patch Volumes: Segmentation-based Consistent Mapping with RGB-D Cameras. In *Third Joint 3DIM/3DPVT Conference (3DV)*.
- Henry, P., Krainin, M., Herbst, E., Ren, X., and Fox, D. (2012). RGB-D mapping: Using Kinect-style depth cameras for dense 3D modeling of indoor environments. *The Int. Journal of Robotics Research*.
- Hornung, A., Wurm, K. M., Bennewitz, M., Stachniss, C., and Burgard, W. (2013). OctoMap: An efficient probabilistic 3D mapping framework based on octrees. *Autonomous Robots*. Software available at <http://octomap.github.com>.
- Hu, G., Huang, S., Zhao, L., Alempijevic, A., and Dissanayake, G. (2012). A robust RGB-D SLAM algorithm. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 1714–1719.
- Huang, A. S., Bachrach, A., Henry, P., Krainin, M., Maturana, D., Fox, D., and Roy, N. (2011). Visual odometry and mapping for autonomous flight using an RGB-D camera. In *Int. Symposium on Robotics Research (ISRR)*, Flagstaff, Arizona, USA.
- Jacobson, A. and Sorkine, O. (2011). Stretchable and twistable bones for skeletal shape deformation. *ACM Transactions on Graphics (proceedings of ACM SIGGRAPH ASIA)*, 30(6):165:1–165:8.
- Kaess, M., Ranganathan, A., and Dellaert, F. (2008). iSAM: Incremental smoothing and mapping. *IEEE Transactions on Robotics (TRO)*, 24(6):1365–1378.
- Karan, K. S. (2000). Skinning characters using surface-oriented free-form deformations. In *In Graphics Interface 2000*, pages 35–42.
- Karpathy, A., Miller, S., and Fei-Fei, L. (2013). Object discovery in 3d scenes via shape analysis. In *International Conference on Robotics and Automation (ICRA)*.
- 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 *Third Joint 3DIM/3DPVT Conference (3DV)*.
- Kerl, C., Sturm, J., and Cremers, D. (2013). Robust odometry estimation for RGB-D cameras. In *Proc. of the IEEE Int. Conf. on Robotics and Automation (ICRA)*.
- Khoshelham, K. and Elberink, S. O. (2012). Accuracy and resolution of kinect depth data for indoor mapping applications. *Sensors*, 12(2):1437–1454.
- Klein, G. and Murray, D. (2007). Parallel tracking and mapping for small AR workspaces. In *Proc. Sixth IEEE and ACM International Symposium on Mixed and Augmented Reality (ISMAR'07)*, Nara, Japan.
- Lee, D., Kim, H., and Myung, H. (2012). GPU-based real-time RGB-D 3D SLAM. In *Ubiquitous Robots and Ambient Intelligence (URAI), 2012 9th International Conference on*, pages 46–48.
- Marton, Z. C., Rusu, R. B., and Beetz, M. (2009). On fast surface

- reconstruction methods for large and noisy datasets. In *Proc. of the IEEE Int. Conf. on Robotics and Automation (ICRA)*, Kobe, Japan.
- Meilland, M. and Comport, A. (2013). On unifying key-frame and voxel-based dense visual SLAM at large scales. In *International Conference on Intelligent Robots and Systems*, Tokyo, Japan. IEEE/RSJ.
- Meilland, M., Drummond, T., and Comport, A. (2013). A Unified Rolling Shutter and Motion Model for Dense 3D Visual Tracking. In *International Conference on Computer Vision*, Sydney, Australia.
- Newcombe, R. A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A. J., Kohli, P., Shotton, J., Hodges, S., and Fitzgibbon, A. (2011). KinectFusion: Real-time Dense Surface Mapping and Tracking. In *Proc. of the 2011 10th IEEE Int. Symposium on Mixed and Augmented Reality, ISMAR '11*, pages 127–136, Washington, DC, USA.
- Nießner, M., Zollhöfer, M., Izadi, S., and Stamminger, M. (2013). Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (TOG)*.
- Pirker, K., Rütger, M., Schweighofer, G., and Bischof, H. (2011). GPSlam: Marrying sparse geometric and dense probabilistic visual mapping. In *Proc. of the British Machine Vision Conf.*, pages 115.1–115.12.
- Roth, H. and Vona, M. (2012). Moving volume KinectFusion. In *British Machine Vision Conf. (BMVC)*, Surrey, UK.
- Salas-Moreno, R. F., Newcombe, R. A., Strasdat, H., Kelly, P. H. J., and Davison, A. J. (2013). SLAM++: Simultaneous localisation and mapping at the level of objects. In *Proc. Computer Vision and Pattern Recognition (CVPR)*.
- Steinbrücker, F., Kerl, C., Sturm, J., and Cremers, D. (2013). Large-scale multi-resolution surface reconstruction from rgb-d sequences. In *IEEE International Conference on Computer Vision (ICCV)*, Sydney, Australia.
- Steinbrücker, F., Sturm, J., and Cremers, D. (2011). Real-Time Visual Odometry from Dense RGB-D Images. In *Workshop on Live Dense Reconstruction with Moving Cameras at the Int. Conf. on Computer Vision (ICCV)*.
- Stückler, J. and Behnke, S. (2013). Multi-resolution surfel maps for efficient dense 3d modeling and tracking. In *Journal of Visual Communication and Image Representation*.
- 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 Int. Conf. on Intelligent Robot Systems (IROS)*.
- Sumner, R. W., Schmid, J., and Pauly, M. (2007). Embedded deformation for shape manipulation. In *ACM SIGGRAPH 2007 papers, SIGGRAPH '07*, New York, NY, USA. ACM.
- Tykkälä, T., Audras, C., and Comport, A. (2011). Direct Iterative Closest Point for Real-time Visual Odometry. In *The Second international Workshop on Computer Vision in Vehicle Technology: From Earth to Mars in conjunction with the International Conference on Computer Vision*, Barcelona, Spain.
- Tykkälä, T., Comport, A. I., and Kamarainen, J.-K. (2013). Photorealistic 3D Mapping of Indoors by RGB-D Scanning Process. In *International Conference on Intelligent Robots and Systems*, Tokyo, Japan.
- van den Braak, G.-J., Nugteren, C., Mesman, B., and Corporaal, H. (2011). Fast hough transform on GPU: Exploration of algorithm trade-offs. In *Advances Concepts for Intelligent Vision Systems*, volume 6915 of *Lecture Notes in Computer Science*, pages 611–622.
- Wagner, R., Frese, U., and Bäuml, B. (2013). 3D Modeling, Distance and Gradient Computation for Motion Planning: A Direct GPGPU Approach. In *International Conference on Robotics and Automation (ICRA)*.
- Whelan, T., Johannsson, H., Kaess, M., Leonard, J., and McDonald, J. (2013a). Robust real-time visual odometry for dense RGB-D mapping. In *IEEE Intl. Conf. on Robotics and Automation, ICRA*, Karlsruhe, Germany.
- Whelan, T., Kaess, M., Fallon, M., Johannsson, H., Leonard, J., and McDonald, J. (2012). Kintinuous: Spatially extended KinectFusion. In *RSS Workshop on RGB-D: Advanced Reasoning with Depth Cameras*, Sydney, Australia.
- Whelan, T., Kaess, M., Leonard, J., and McDonald, J. (2013b). Deformation-based loop closure for large scale dense RGB-D SLAM. In *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems, IROS*, Tokyo, Japan.
- Zeng, M., Zhao, F., Zheng, J., and Liu, X. (2012). A Memory-Efficient KinectFusion Using Octree. In *Computational Visual Media*, volume 7633 of *Lecture Notes in Computer Science*, pages 234–241. Springer.
- Zhou, Q. and Koltun, V. (2013). Dense scene reconstruction with points of interest. In *SIGGRAPH 2013*, Anaheim, CA, USA. ACM.
- Zhou, Q., Miller, S., and Koltun, V. (2013). Elastic Fragments for Dense Scene Reconstruction. In *Int. Conf. on Computer Vision (ICCV)*.

## A Index to Multimedia Extensions

The multimedia extensions to this article are at: <http://www.ijrr.org>.

Extension	Type	Description
1	Video	Indoors and Two floors dataset deformation visualisations.
2	Video	In/outdoors dataset full real-time reconstruction.
3	Video	Apartment dataset full real-time reconstruction.

## B Algorithms

---

**Algorithm 1:** Color Integration

---

**Input:**  $\mathbf{rgb}_i$  Current RGB image  
 $\mathbf{d}_i$  Current depth map  
 $\mathbf{n}_i$  Current normal map  
 $S(\mathbf{s})_i$  Current TSDF volume  
 $\mathbf{s} \in \Psi$  Current voxel  
 $\mathbf{p} \in \Omega$  Current pixel

**do**  
   $c \leftarrow 0$   
  **for each**  $\mathbf{p}_k$  in  $7 \times 7$  area around  $\mathbf{p}$  **do**  
    **if**  $|\mathbf{d}_i(\mathbf{p}_k) - \mathbf{d}_i(\mathbf{p})| > \text{depth threshold}$  **or**  $\mathbf{d}_i(\mathbf{p}_k) = 0$  **then**  
       $c \leftarrow c + 1$   
  **if**  $c < \text{count threshold}$  **then**  
     $w_c = \min(1.0, \mathbf{n}_i(\mathbf{p})_z / \text{max\_weight})$   
     $S(\mathbf{s})_i^{R'} = (S(\mathbf{s})_{i-1}^W S(\mathbf{s})_{i-1}^R + w_c \mathbf{rgb}_i(\mathbf{p})^R) / (S(\mathbf{s})_{i-1}^W + w_c)$   
     $S(\mathbf{s})_i^{G'} = (S(\mathbf{s})_{i-1}^W S(\mathbf{s})_{i-1}^G + w_c \mathbf{rgb}_i(\mathbf{p})^G) / (S(\mathbf{s})_{i-1}^W + w_c)$   
     $S(\mathbf{s})_i^{B'} = (S(\mathbf{s})_{i-1}^W S(\mathbf{s})_{i-1}^B + w_c \mathbf{rgb}_i(\mathbf{p})^B) / (S(\mathbf{s})_{i-1}^W + w_c)$   
**end**

---

---

**Algorithm 2:** Interest Point Accumulation

---

**Input:**  $\frac{\partial I_n}{\partial x}$  and  $\frac{\partial I_n}{\partial y}$  intensity image derivatives  
 $s$  minimum gradient scale for pyramid level

**Output:**  $\mathcal{L}$  list of interest points  
 $k_{\mathcal{L}}$  global point count

**Data:**  $\alpha$  thread block x-dimension  
 $\beta$  thread block y-dimension  
 $\gamma$  pixels per thread  
 $\iota$  shared memory local list  
 $\kappa$  shared memory local index  
 $\text{blockIdx}$  CUDA block index  
 $\text{threadIdx}$  CUDA thread index

**in parallel do**  
   $i \leftarrow \beta * \text{blockIdx.y} + \text{threadIdx.y}$   
   $j \leftarrow \alpha * \gamma * \text{blockIdx.x} + \gamma * \text{threadIdx.x}$   
  **if**  $\text{threadIdx.x} = 0$  **and**  $\text{threadIdx.y} = 0$  **then**  
     $\kappa \leftarrow 0$   
  **syncthreads()**  
  **for**  $l \leftarrow 0$  **to**  $\gamma$  **do**  
     $\mathbf{p} \leftarrow (i, j + l)$   
     $g^2 = \frac{\partial I_n}{\partial x}(\mathbf{p})^2 + \frac{\partial I_n}{\partial y}(\mathbf{p})^2$   
    **if**  $g^2 \geq s$  **then**  
       $\text{idx} \leftarrow \text{atomicInc}(\kappa)$   
       $\iota_{\text{idx}} \leftarrow \mathbf{p}$   
  **syncthreads()**  
   $b \leftarrow \alpha * \gamma * \text{threadIdx.y} + \gamma * \text{threadIdx.x}$   
  **for**  $l \leftarrow 0$  **to**  $\gamma$  **do**  
     $a \leftarrow b + l$   
    **if**  $a < \kappa$  **then**  
       $\text{idx} \leftarrow \text{atomicInc}(k_{\mathcal{L}})$   
       $\mathcal{L}_{\text{idx}} \leftarrow \iota_a$   
**end**

---

---

**Algorithm 3:** Correspondence Accumulation

---

**Input:**  $\mathcal{L}$  list of interest points  
 $d_\delta$  maximum change in point depth  
 $[I_{n-1}, M_{n-1}]$  previous intensity depth pair  
 $[I_n, M_n]$  current intensity depth pair  
 $\mathbf{R}^l$  camera rotation in image  
 $\mathbf{t}^l$  camera translation in image

**Output:**  $C$  correspondence list of the form  $(\mathbf{p}, \mathbf{p}', \Delta)$   
 $k_C$  global point count  
 $\sigma$  global intensity difference sum

**Data:**  $\alpha$  thread block x-dimension  
 $\gamma$  pixels per thread  
 $\iota$  shared memory local list  
 $\kappa$  shared memory local index  
 $\text{blockIdx}$  CUDA block index  
 $\text{threadIdx}$  CUDA thread index

**in parallel do**  
   $i \leftarrow \alpha * \gamma * \text{blockIdx.x} + \gamma * \text{threadIdx.x}$   
  **if**  $\text{threadIdx.x} = 0$  **then**  
     $\kappa \leftarrow 0$   
  **syncthreads()**  
  **for**  $l \leftarrow 0$  **to**  $\gamma$  **do**  
     $\mathbf{p} \leftarrow \mathcal{L}_{i+l}$   
     $z \leftarrow M_n(\mathbf{p})$   
    **if**  $\text{isValid}(z)$  **then**  
       $(x', y', z')^T \leftarrow z(\mathbf{R}^l(\mathbf{p}, 1)^T) + \mathbf{t}^l$   
       $\mathbf{p}' \leftarrow (\frac{x'}{z'}, \frac{y'}{z'})^T$   
      **if**  $\text{isInImage}(\mathbf{p}')$  **then**  
         $d \leftarrow M_{n-1}(\mathbf{p}')$   
        **if**  $\text{isValid}(d)$  **and**  $|z' - d| \leq d_\delta$  **then**  
           $\text{idx} \leftarrow \text{atomicInc}(\kappa)$   
           $\iota_{\text{idx}} \leftarrow (\mathbf{p}, \mathbf{p}', I_n(\mathbf{p}) - I_{n-1}(\mathbf{p}'))$   
  **syncthreads()**  
   $b \leftarrow \gamma * \text{threadIdx.x}$   
  **for**  $l \leftarrow 0$  **to**  $\gamma$  **do**  
     $a \leftarrow b + l$   
    **if**  $a < \kappa$  **then**  
       $\text{atomicAdd}(\sigma, \iota_a^2)$   
       $\text{idx} \leftarrow \text{atomicInc}(k_C)$   
       $C_{\text{idx}} \leftarrow \iota_a$   
**end**

---

---

**Algorithm 4:** Incremental Deformation Node Sampling

---

**Input:**  $P$  camera pose graph made up of  $\mathbf{R}_i$  and  $\mathbf{t}_i$   
 $i$  pose id of last added node  
 $d_p$  pose sampling rate

**Output:**  $N$  set of deformation graph nodes

**do**  
   $l \leftarrow |N|$   
  **if**  $l = 0$  **then**  
     $N_l^g \leftarrow \mathbf{t}_0$   
     $l \leftarrow l + 1$   
     $i \leftarrow 0$   
   $P_{\text{last}} \leftarrow P_i$   
  **for**  $i$  **to**  $|P|$  **do**  
    **if**  $\|\mathbf{t}_i - \mathbf{t}_{\text{last}}\|_2 > d_p$  **then**  
       $N_l^g \leftarrow \mathbf{t}_i$   
       $l \leftarrow l + 1$   
       $P_{\text{last}} \leftarrow P_i$   
**end**

---

---

**Algorithm 5:** Back-Traversal Vertex Association

---

**Input:**  $C$  cloud slices

$N$  set of deformation graph nodes

$b_p$  number of poses to traverse back

$P_{C_j}$  pose associated with cloud slice  $C_j$

**Output:**  $\mathcal{N}(v)$  for each  $v$

**do**

**foreach**  $C_j$  **do**

**foreach**  $v \in C_j$  **do**

$l \leftarrow \text{binary\_search\_closest}(P_{C_j}, N)$

$N' \leftarrow \emptyset$

$n \leftarrow 0$

**for**  $i \leftarrow 0$  **to**  $b_p$  **do**

$N'_n \leftarrow N_l$

$n \leftarrow n + 1$

$l \leftarrow l - 1$

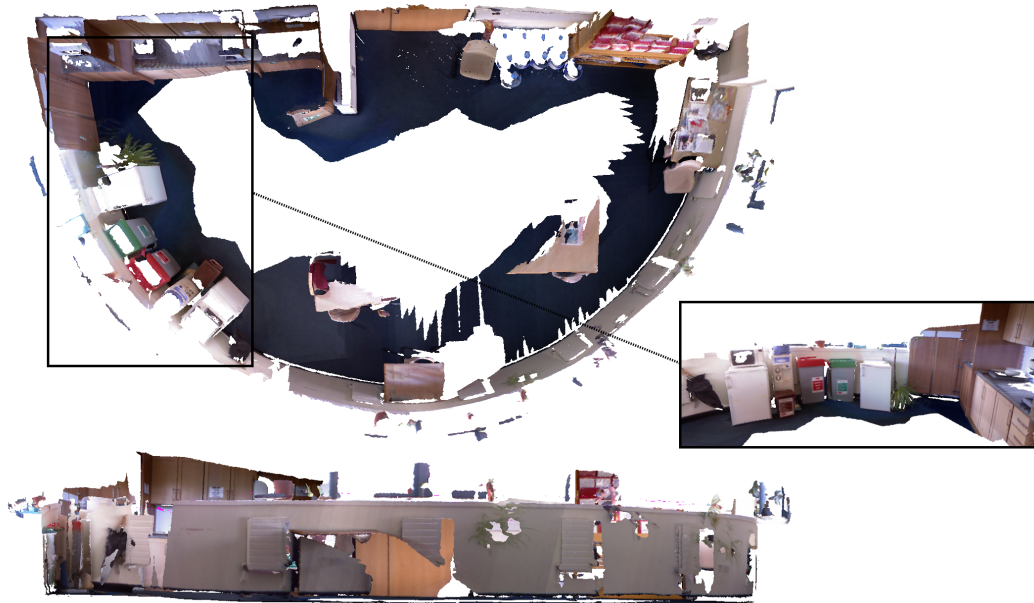
$\text{sort\_by\_distance}(N', v)$

$\mathcal{N}(v) \leftarrow N'_{1 \rightarrow k}$

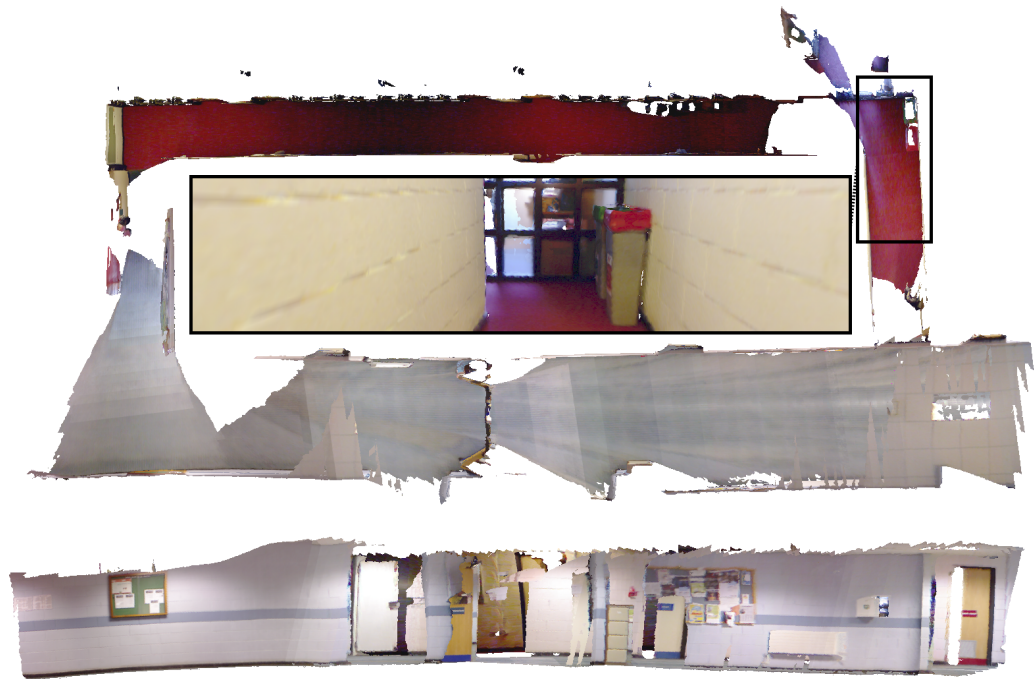
**end**

---

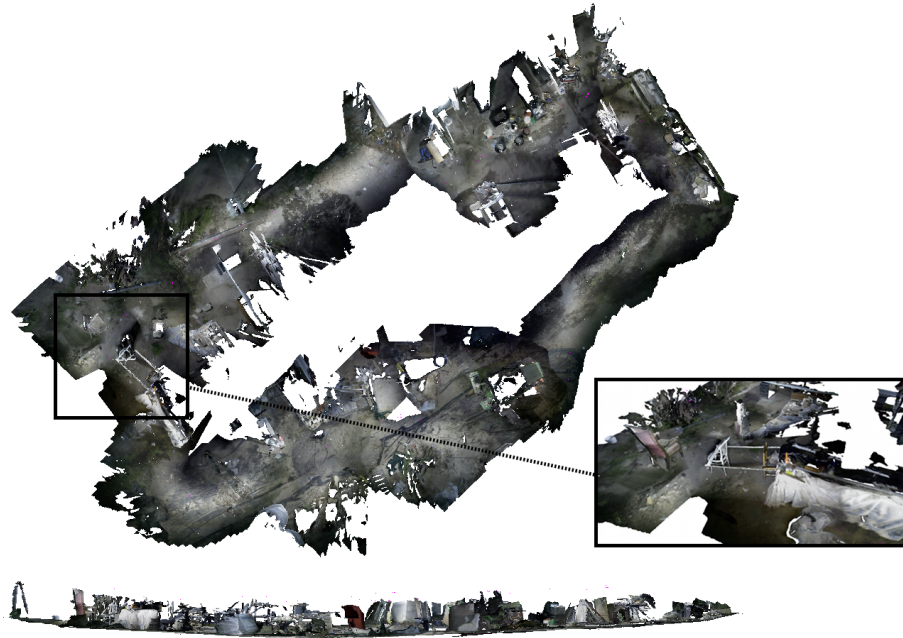




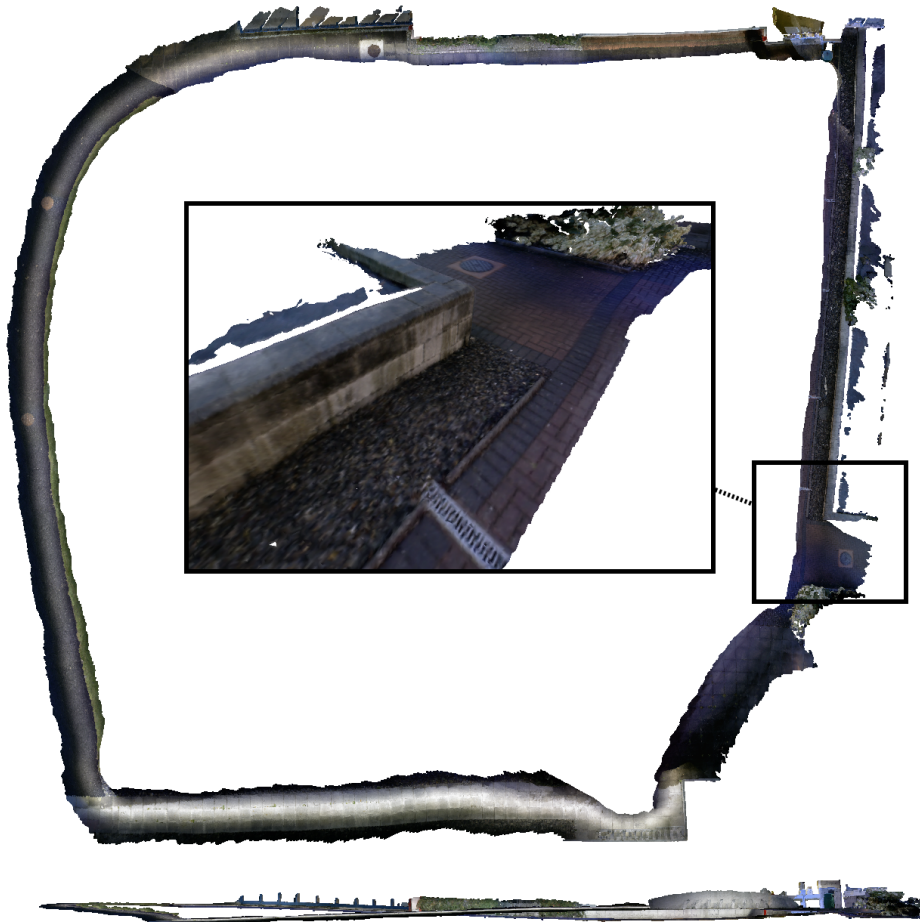
**Figure 22:** Dataset of a small coffee room. Inset shows everyday objects such as bins and fridges are captured in high detail and how the deformation approach works well in smaller environments.



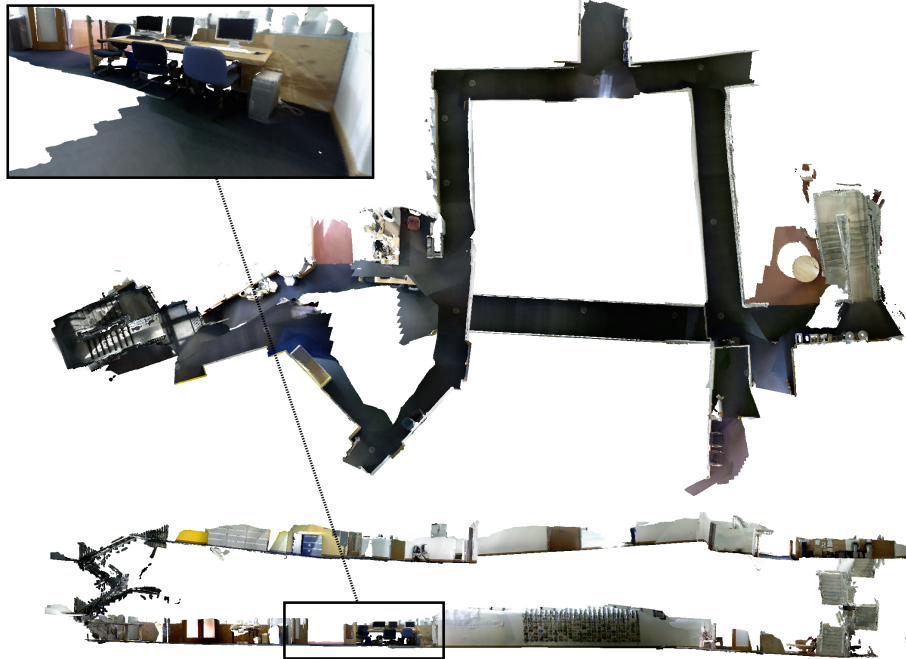
**Figure 23:** Corridor loop closure dataset. The inset shows map consistency at the point of loop closure. Multimedia Extension 1 shows the actual map correcting deformation occurring.



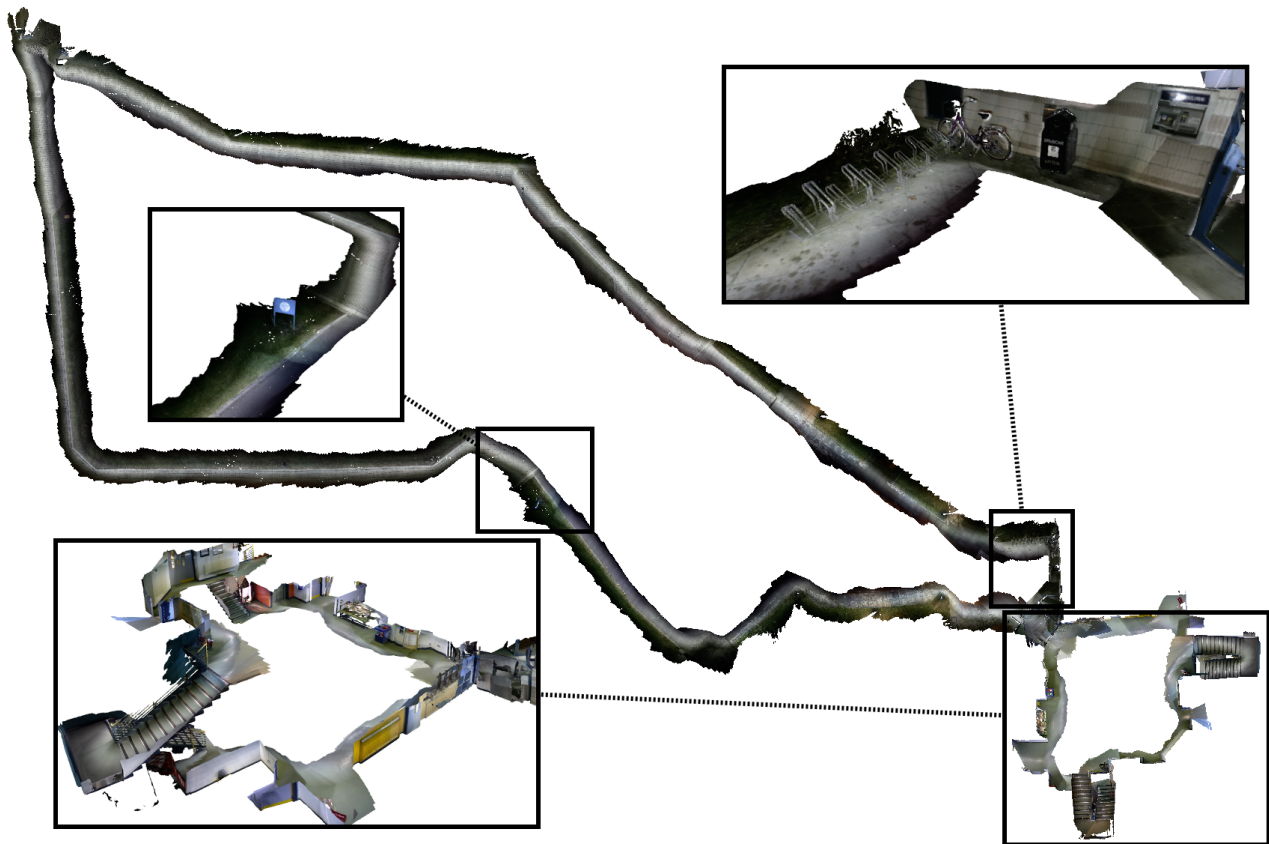
**Figure 24:** Large cluttered outdoor dataset. Inset shows chairs and metal bars are reconstructed well.



**Figure 25:** Large outdoor dataset. Inset shows brickwork is clearly visible.



**Figure 26:** Dataset composed of two floors. Inset shows everyday objects such as chairs and computers are captured in high detail. Multimedia Extension 1 shows the actual map correcting deformation occurring.



**Figure 27:** Large indoor and outdoor dataset made up of over five million vertices. Insets show the high fidelity of small scale features in the map. Multimedia Extension 2 shows this entire dataset running from start to finish in real-time, including online loop closure.



**Figure 28:** Sequence over two floors of an apartment with over six million vertices. Small details such as bathroom fixtures and objects around the environment are clearly reconstructed. Multimedia Extension 3 shows this dataset running from start to finish in real-time, including online loop closure.