

BundleFusion: Real-time Globally Consistent 3D Reconstruction using On-the-fly Surface Re-integration

ANGELA DAI and MATTHIAS NIEßNER

Stanford University

MICHAEL ZOLLHÖFER

Max-Planck-Institute for Informatics

SHAHRAM IZADI

Microsoft Research

and

CHRISTIAN THEOBALT

Max-Planck-Institute for Informatics

Real-time, high-quality, 3D scanning of large-scale scenes is key to mixed reality and robotic applications. However, scalability brings challenges of drift in pose estimation, introducing significant errors in the accumulated model. Approaches often require hours of offline processing to globally correct model errors. Recent online methods demonstrate compelling results, but suffer from: (1) needing minutes to perform online correction preventing true real-time use; (2) brittle frame-to-frame (or frame-to-model) pose estimation resulting in many tracking failures; or (3) supporting only unstructured point-based representations, which limit scan quality and applicability. We systematically address these issues with a novel, real-time, end-to-end reconstruction framework. At its core is a robust pose estimation strategy, optimizing per frame for a global set of camera poses by considering the complete history of RGB-D input with an efficient hierarchical approach. We remove the heavy reliance on temporal tracking, and continually localize to the globally optimized frames instead. We contribute a parallelizable optimization framework, which employs correspondences based on sparse features and dense geometric and photometric matching. Our approach estimates globally optimized (i.e., bundle adjusted) poses in real-time, supports robust tracking with recovery from gross tracking failures (i.e., relocalization), and re-estimates the 3D model in real-time to ensure global consistency; all within a single framework. Our approach outperforms state-of-the-art online systems with quality on par to offline methods, but with unprecedented speed and scan completeness. Our framework leads to a comprehensive online scanning solution for large indoor environments, enabling ease of use and high-quality results¹.

Categories and Subject Descriptors: I.3.3 [Computer Graphics]: Picture/Image Generation—*Digitizing and Scanning*

Additional Key Words and Phrases: RGB-D, scan, real-time, global consistency, scalable

1. INTRODUCTION

We are seeing a renaissance in 3D scanning, fueled both by applications such as fabrication, augmented and virtual reality, gaming and robotics, and by the ubiquity of RGB-D cameras, now even available in consumer-grade mobile devices. This has opened up the need for *real-time* scanning at *scale*. Here, the user or robot must scan an entire room (or several spaces) in real-time, with instantaneous and

continual integration of the accumulated 3D model into the desired application, whether that is robot navigation, mapping the physical world into the virtual, or providing immediate user feedback during scanning.

However, despite the plethora of reconstruction systems, we have yet to see a single holistic solution for the problem of real-time 3D reconstruction at scale that makes scanning easily accessible to untrained users. This is due to the many requirements that such a solution needs to fulfill:

High-quality surface modeling. We need a single textured and noise-free 3D model of the scene, consumable by standard graphics applications. This requires a high-quality representation that can model *continuous surfaces* rather than discrete points.

Scalability. For mixed reality and robot navigation scenarios, we need to acquire models of entire rooms or several large spaces. Our underlying representation therefore must handle both small- and large-scale scanning while preserving both global structure and maintaining high local accuracy.

Global model consistency. With scale comes the need to correct pose estimation errors and drift, and the subsequent distortions in the acquired 3D model. This correction is particularly challenging at real-time rates, but is key for allowing online revisiting of previously scanned areas or loop closure during actual use.

Robust camera tracking. Apart from incremental errors, camera tracking can also fail in featureless regions. In order to recover, we require the ability to relocalize. Many existing approaches rely heavily on proximity to the previous frame, limiting fast camera motion and recovery from tracking failure. Instead, we need to (re)localize in a robust manner without relying on temporal coherence.

On-the-fly model updates. In addition to robust tracking, input data needs to be integrated to a 3D representation and interactively visualized. The challenge is to update the model after data has been integrated, in accordance with new pose estimates.

Real-time rates. The ability to react to instantaneous feedback is crucial to 3D scanning and key to obtaining high-quality results. The real-time capability of a 3D scanning method is fundamental to AR/VR and robotics applications.

Researchers have studied specific parts of this problem, but to date there is no single approach to tackle all of these requirements in real time. This is the very aim of this paper, to systematically address *all* these requirements in a single, end-to-end real-time reconstruction framework. At the core of our method is a robust pose estimation strategy, which globally optimizes for the camera

¹Our source code and all reconstruction results are publicly available: <http://graphics.stanford.edu/projects/bundlefusion/>

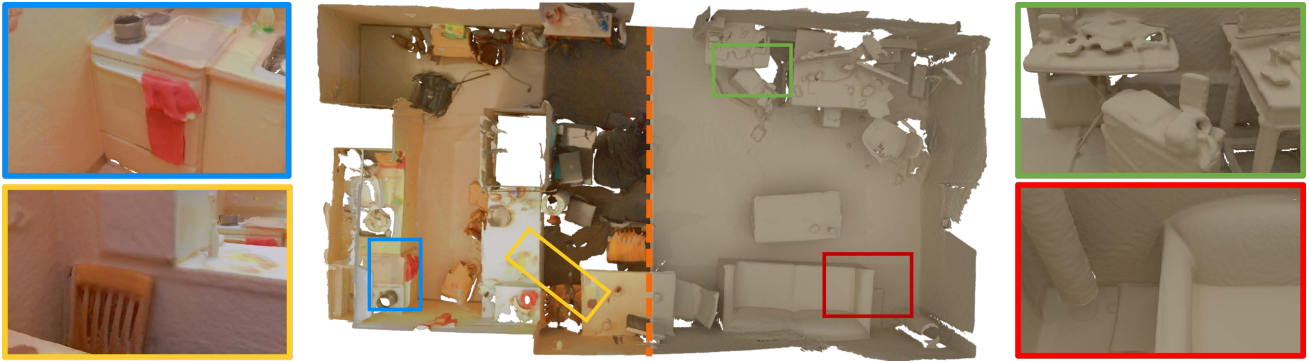


Fig. 1. Our novel real-time 3D reconstruction approach solves for global pose alignment and obtains dense volumetric reconstructions at a level of quality and completeness that was previously only attainable with offline approaches.

trajectory per frame, considering the complete history of RGB-D input in an efficient *local-to-global* hierarchical optimization framework. Since we globally correlate each RGB-D frame, loop closure is handled implicitly and continuously, removing the need for any explicit loop closure detection. This enables our method to be extremely robust to tracking failures, with tracking far less brittle than existing frame-to-frame or frame-to-model RGB-D approaches. If tracking failures occur, our framework instantaneously relocalizes in a globally consistent manner, even when scanning is interrupted and restarted from a completely different viewpoint. Areas can also be revisited multiple times without problem, and reconstruction quality continuously improves. This allows for a robust scanning experience, where even novice users can perform large-scale scans without failure.

Key to our work is a new fully parallelizable *sparse-then-dense* global pose optimization framework: sparse RGB features are used for coarse global pose estimation, ensuring proposals fall within the basin of convergence of the following dense step which considers both photometric and geometric consistency for fine-scale alignment. Thus, we maintain global structure with implicit loop closures while achieving high local reconstruction accuracy. To achieve the corresponding model correction, we extend a scalable variant of real-time volumetric fusion [Nießner et al. 2013], but importantly support model updates based on refined poses from our global optimization. Thus, we can correct errors in the 3D model in real time and revisit existing scanned areas. We demonstrate how our approach outperforms current state-of-the-art online systems at unprecedented speed and scan completeness, and even surpasses the accuracy and robustness of offline methods in many scenarios. This leads to a comprehensive real-time scanning solution for large indoor environments, that requires little expertise to operate, making 3D scanning easily accessible to the masses.

In summary, the main contributions of our work are as follows:

(1) A novel, real-time global pose alignment framework which considers the complete history of input frames, removing the brittle and imprecise nature of temporal tracking approaches, while achieving scalability by a rapid hierarchical decomposition of the problem by using a *local-to-global optimization* strategy.

(2) A *sparse-to-dense alignment* strategy enabling both consistent global structure with implicit loop closures and highly-accurate fine-scale pose alignment to facilitate local surface detail.

(3) A new RGB-D re-integration strategy to enable on-the-fly and continuous 3D model updates when refined global pose estimates are available.

(4) Large-scale reconstruction of geometry and texture, demonstrating model refinement in revisited areas, recovery from tracking failures, and robustness to drift and continuous loop closures.

2. RELATED WORK

There has been extensive work on 3D reconstruction over the past decades. Key to high-quality 3D reconstruction is the choice of underlying representation for fusing multiple sensor measurements. Approaches range from unstructured point-based representations [Rusinkiewicz et al. 2002; Weise et al. 2009; Henry et al. 2010; Keller et al. 2013; Whelan et al. 2015], 2.5D depth map [Merrell et al. 2007; Meilland et al. 2013] or height-field [Gallup et al. 2010] methods, to volumetric approaches, based on occupancy grids [Elfes and Matthies 1987; Wurm et al. 2010] or implicit surfaces [Hilton et al. 1996; Curless and Levoy 1996]. While each has trade-offs, volumetric methods based on implicit truncated signed distance fields (TSDFs) have become the de facto method for highest quality reconstructions; e.g., [Levoy et al. 2000; Fuhrmann and Gesele 2014; Fioraio et al. 2015]. They model continuous surfaces, systematically regularize noise, remove the need for explicit topology bookkeeping, and efficiently perform incremental updates. The most prominent recent example is KinectFusion [Newcombe et al. 2011; Izadi et al. 2011] where real-time volumetric fusion of smaller scenes was demonstrated.

One inherent issue with these implicit volumetric methods is their lack of scalability due to reliance on a uniform grid. This has become a focus of much recent research [Roth and Vona 2012; Whelan et al. 2012; Zeng et al. 2012; Chen et al. 2013; Keller et al. 2013; Nießner et al. 2013; Steinbruecker et al. 2014; Reichl et al. 2015], where real-time efficient data structures for volumetric fusion have been proposed. These exploit the sparsity in the TSDF representation to create more efficient spatial subdivision strategies. While this allows for volumetric fusion at scale, pose estimates suffer from drift, causing distortions in the 3D model. Even small pose errors, seemingly negligible on a small local scale, can accumulate to dramatic error in the final 3D model [Nießner et al. 2013].

Most of the research on achieving globally consistent 3D models at scale from RGB-D input requires offline processing and access to all input frames. [Li et al. 2013; Zhou and Koltun 2013; Zhou et al. 2013; Zhou and Koltun 2014; Choi et al. 2015] provide for globally consistent models by optimizing across the entire pose trajectory, but require minutes or even hours of processing time, meaning *real-time* revisiting or refinement of reconstructed areas is infeasible.

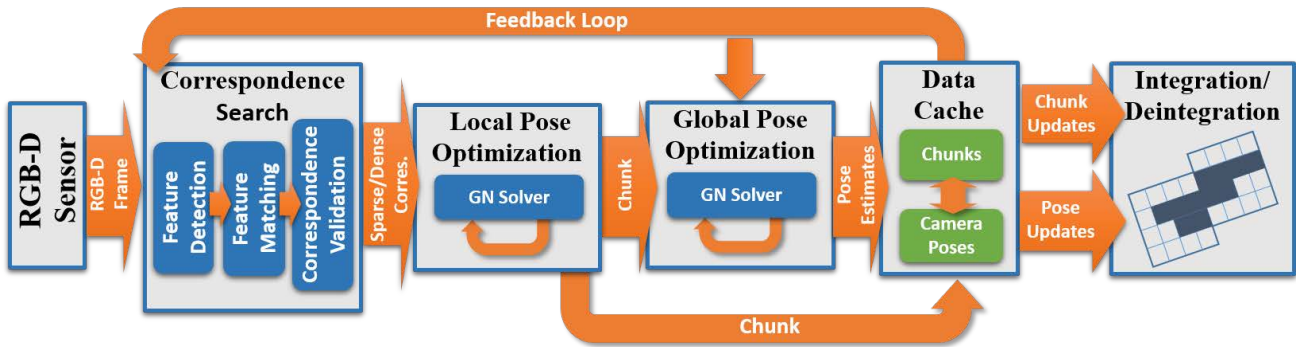


Fig. 2. Our global pose optimization takes as input the RGB-D stream of a commodity sensor, detects pairwise correspondences between the input frames, and performs a combination of local and global alignment steps using sparse and dense correspondences to compute per-frame pose estimates.

Real-time, drift-free pose estimation is a key focus in the simultaneous localization and mapping (SLAM) literature. Many real-time monocular RGB methods have been proposed, including sparse methods [Klein and Murray 2007], semi-dense [Engel et al. 2013; Forster et al. 2014] or direct methods [Meilland et al. 2011; Engel et al. 2014]. Typically these approaches rely on either pose-graph optimization [Kümmerle et al. 2011] or bundle adjustment [Triggs et al. 2000], minimizing reprojection error across frames and/or distributing the error across the graph. While impressive tracking results have been shown using only monocular RGB sensors, these approaches do not generate detailed dense 3D models, which is the aim of our work.

MonoFusion [Pradeep et al. 2013] augments sparse SLAM bundle adjustment with dense volumetric fusion, showing compelling monocular results but on small-scale scenes. Real-time SLAM approaches typically first estimate poses *frame-to-frame* and perform correction in a background thread (running slower than real-time rates; e.g., 1Hz). In contrast, DTAM [Newcombe et al. 2011] uses the concept of *frame-to-model* tracking (from KinectFusion [Newcombe et al. 2011; Izadi et al. 2011]) to estimate the pose directly from the reconstructed dense 3D model. This omits the need for a correction step, but clearly does not scale to larger scenes.

Pose estimation from range data typically is based on variants of the iterative closest point (ICP) algorithm [Besl and McKay 1992; Rusinkiewicz and Levoy 2001]. In practice, this makes tracking extremely brittle and has led researchers to explore either the use of RGB data to improve frame-to-frame tracking [Whelan et al. 2013] or the use of global pose estimation correction, including pose graph optimization [Steinbruecker et al. 2013], loop closure detection [Whelan et al. 2013], incremental bundle adjustment [Whelan et al. 2015; Fioraio et al. 2015], or recovery from tracking failures by image or keypoint-based relocalization [Glocker et al. 2015; Valentin et al. 2015].

These systems are state-of-the-art in terms of online correction of both pose and underlying 3D model. However, they either require many seconds or even minutes to perform online optimization [Whelan et al. 2013; Fioraio et al. 2015]; assume very specific camera trajectories to detect explicit loop closures limiting free-form camera motions and scanning [Whelan et al. 2013]; rely on computing optimized camera poses prior to fusion limiting the ability to refine the model afterwards [Steinbruecker et al. 2013], or use point-based representations that limit model quality and lack general applicability where continuous surfaces are needed [Whelan et al. 2015].

3. METHOD OVERVIEW

The core of our approach is an efficient global pose optimization algorithm which operates in unison with a large-scale, real-time 3D reconstruction framework; see Fig. 2. At every frame, we continuously run pose optimization and update the reconstruction according to the newly-computed pose estimates. We do not strictly rely on temporal coherence, allowing for free-form camera paths, instantaneous relocalization, and frequent revisiting of the same scene region. This makes our approach robust towards sensor occlusion, fast frame-to-frame motions and featureless regions.

We take as input the RGB-D stream captured by a commodity depth sensor. To obtain global alignment, we perform a sparse-then-dense global pose optimization: we use a set of sparse feature correspondences to obtain a coarse global alignment, as sparse features inherently provide for loop closure detection and relocalization. This alignment is then refined by optimizing for dense photometric and geometric consistency. Sparse correspondences are established through pairwise Scale-Invariant Feature Transform (SIFT) [Lowe 2004] feature correspondences between all input frames (see Sec. 4.1). That is, detected SIFT keypoints are matched against all previous frames, and filtered to remove outliers (see Sec. 4.1.1).

To make real-time global pose alignment tractable, we perform a hierarchical local-to-global pose optimization (see Sec. 4.2) using the filtered frame correspondences. On the first hierarchy level, every consecutive n frames compose a *chunk*, which is locally pose optimized under the consideration of its contained frames. On the second hierarchy level, all chunks are correlated with respect to each other and globally optimized. This is akin to hierarchical submapping [Maier et al. 2014]; however, instead of analyzing global connectivity once all frames are available, our new method forms *chunks* based on the current temporal window.

This hierarchical two-stage optimization strategy reduces the number of unknowns per optimization step and ensures our method scales to large scenes. Pose alignment on both levels is formulated as energy minimization problem in which both the filtered sparse correspondences, as well as dense photometric and geometric constraints are considered (see Sec. 4.3). To solve this highly-nonlinear optimization problem on both hierarchy levels, we employ a fast data-parallel GPU-solver tailored to the problem (see Sec. 4.4).

A dense scene reconstruction is obtained using a sparse volumetric representation and fusion [Nießner et al. 2013], which scales to large scenes in real-time. The continuous change in the optimized global pose necessitates continuous updates to the global 3D scene representation (see Sec. 5). A key novelty is to allow for symmetric

on-the-fly reintegration of RGB-D frames. In order to update the pose of a frame with an improved estimate, we remove the RGB-D image at the old pose with a new real-time *de-integration* step, and *reintegrate* the RGB-D image at its new pose. Thus, the volumetric model continuously improves as more RGB-D frames and refined pose estimates become available; e.g., if a loop is closed (cf. Fig. 12).

4. GLOBAL POSE ALIGNMENT

We first describe the details of our real-time global pose optimization strategy, which is the foundation for *online*, globally-consistent 3D reconstruction. Input to our approach is the live RGB-D stream $\mathbf{S} = \{f_i = (\mathcal{C}_i, \mathcal{D}_i)\}_i$ captured by a commodity sensor. We assume spatially and temporally aligned color \mathcal{C}_i and depth data \mathcal{D}_i in each frame, captured at 30Hz and 640×480 pixel resolution. The goal is to find a set of 3D correspondences between the frames in the input sequence, and then to find an optimal set of rigid camera transforms $\{\mathcal{T}_i\}$ such that all frames align as best as possible. The transformation $\mathcal{T}_i(\mathbf{p}) = \mathbf{R}_i\mathbf{p} + \mathbf{t}_i$ (rotation \mathbf{R}_i , translation \mathbf{t}_i) maps from the local camera coordinates of the i -th frame to the world space coordinate system; we assume the first frame defines the world coordinate system.

4.1 Feature Correspondence Search

In our framework, we first search for sparse correspondences between frames using efficient feature detection, feature matching, and correspondence filtering steps. These sparse correspondences are later used in tandem with dense photometric correspondences, but since accurate sparse correspondences are crucial to attaining the basin of convergence of the dense optimization, we elaborate on their search and filtering below. For each new frame, SIFT features are detected and matched to the features of all previously seen frames. We use SIFT as it accounts for the major variation encountered during hand-held RGB-D scanning, namely: image translation, scaling, and rotation. Matches between each pair of frames are then filtered to produce a list of valid pairwise correspondences as input to global pose optimization. Our correspondence search is performed entirely on the GPU, avoiding the overhead of copying data (e.g., feature locations, descriptors, matches) to the host. We compute SIFT keypoints and descriptors at 4 – 5 ms per frame, and match a pair of frames in ≈ 0.05 ms (in parallel). We can thus find full correspondences in real-time against up to over 20K frames, matched in a hierarchical fashion, for every new input RGB-D image.

4.1.1 Correspondence Filtering. To minimize outliers, we filter the sets of detected pairwise correspondences based on geometric and photometric consistency. Note that further robustness checks are built into the optimization (not described in this section; see Sec. 4.4.1 for details).

Key Point Correspondence Filter For a pair of frames f_i and f_j with detected corresponding 3D points P from f_i , and Q from f_j , the key point correspondence filter finds a set of correspondences which exhibit a stable distribution and a consistent rigid transform. Correspondences are greedily aggregated (in order of match distance); for each newly added correspondence, we compute the rigid transform $\mathcal{T}_{ij}(\mathbf{p}) = (\mathcal{T}_j^{-1} \circ \mathcal{T}_i)(\mathbf{p})$, which minimizes the RMSD between the current set of correspondences P_{cur} and Q_{cur} , using the Kabsch algorithm [Kabsch 1976; Gower 1975]. We further check whether this is an ambiguously determined transform (e.g. the correspondences lie on a line or exhibit rotational symmetry) by performing a condition analysis of the distribution of points of P_{cur}

and Q_{cur} as well as the cross-covariance between P_{cur} and Q_{cur} ; if any of these condition numbers are high then the system is considered unstable. Thus, if the re-projection error under \mathcal{T}_{ij} is high or the condition analysis determines instability, then correspondences are removed (in order of re-projection error) until this is not the case anymore or there are too few correspondences to determine a rigid transform. If the resulting set of correspondences for f_i and f_j do not produce a valid transform, all correspondences between f_i and f_j are discarded.

Surface Area Filter In addition, we check that the surface spanned by the features is large enough, as correspondences spanning small physical size are prone to ambiguity. For frames f_i and f_j , we estimate the surface areas spanned by the 3D keypoints P of f_i and by the 3D keypoints Q of f_j . For each set of 3D points, we project them into the plane given by their two principal axes, with surface area given by the 2D oriented bounding box of the resulting projected points. If the areas spanned by P and Q are insufficient ($< 0.032\text{m}^2$), the set of matches is deemed ambiguous and discarded.

Dense Verification Finally, we perform a dense two-sided geometric and photometric verification step. For frames f_i and f_j , we use the computed relative transform \mathcal{T}_{ij} from the key point correspondence filter to align the coordinate systems of f_i and f_j . We then measure the average depth discrepancy, normal deviation and photoconsistency of the re-projection in both directions. For efficiency reasons, this step is performed on filtered and downsampled input frames of size $w' \times h' = 80 \times 60$. Note that when a new RGB-D image f_i arrives, its filtered and downsampled color intensity \mathcal{C}_i^{low} and depth \mathcal{D}_i^{low} are cached for efficiency. The camera space positions P_i^{low} and normals N_i^{low} of each \mathcal{D}_i^{low} are also computed and cached per frame. With π denoting the camera intrinsics for the downsampled images, the total re-projection error from f_i to f_j is:

$$E_r(f_i, f_j) = \sum_{x,y} \left\| \mathcal{T}_{ij}(\mathbf{p}_{i,x,y}) - \mathbf{q}_{j,x,y} \right\|_2.$$

Here, $\mathbf{p}_{i,x,y} = P_i^{low}(x, y)$ and $\mathbf{q}_{j,x,y} = P_j^{low}(\pi^{-1}(\mathcal{T}_{ij}\mathbf{p}_{i,x,y}))$. Dense correspondences are determined to be invalid when their positions, normals, or colors do not correspond (i.e., accounting for possible occlusion error). Matches between f_i and f_j are invalidated in the case of excessive re-projection error ($> 0.075\text{m}$) or insufficient valid correspondences ($< 0.02w'h'$). This check is efficiently implemented with a single kernel call, such that each thread block handles one pair of images, with re-projection computed through local reductions.

If all checks are passed, the correspondences are added to the valid set, which is used later on for pose optimization. We only consider a frame-to-frame match if the valid set comprises at least N_{min} correspondences between them. Note that $N_{min} = 3$ is sufficient to define a valid frame-to-frame transform; however, we found $N_{min} = 5$ to be a good compromise between precision and recall.

4.2 Hierarchical Optimization

In order to run at real-time rates on up to tens of thousands of RGB-D input frames, we apply a hierarchical optimization strategy. The input sequence is split into short *chunks* of consecutive frames. On the lowest hierarchy level, we optimize for local alignments within a *chunk*. On the second hierarchy level, chunks are globally aligned against each other, using representative *keyframes* with associated features per chunk.

Local Intra-Chunk Pose Optimization Intra-chunk alignment is based on chunks of $N_{chunk} = 11$ consecutive frames in the input

RGB-D stream; adjacent chunks overlap by 1 frame. The goal of local pose optimization is to compute the best intra-chunk alignments $\{\mathcal{T}_i\}$, relative to the first frame of the chunk, which locally defines the reference frame. To this end, valid feature correspondences are searched between all pairs of frames of the chunk, and then the energy minimization approach described in Sec. 4.3 is applied, jointly considering both these feature correspondences *and* dense photometric and geometric matching. Since each chunk only contains a small number of consecutive frames, the pose variation within the chunk is small, and we can initialize each of the \mathcal{T}_i to the identity matrix. To ensure that the local pose optimization result after convergence is sufficiently accurate, we apply the Dense Verification test (see Sec. 4.1.1) to each pair of images within the chunk using the optimized local trajectory. If the re-projection error is too large for any pair of images ($> 0.05\text{m}$), the chunk is discarded and not used in the global optimization.

Per-Chunk Keyframes Once a chunk has been completely processed, we define the RGB-D data from the first frame in the chunk to be the chunk’s *keyframe*. We also compute a representative aggregate *keyframe feature set*. Based on the optimized pose trajectory of the chunk, we compute a coherent set of 3D positions of the inter-chunk feature points in world space. These 3D positions may contain multiple instances of the same real-world point, found in separate pairwise frame matches. Thus to obtain the keyframe feature set, multiple feature point instances which match in their feature descriptors and coincide in 3D world space are merged to one best 3D representative in the least squares sense. This keyframe feature set is mapped into the space of the keyframe using the respective transformation. Note that once this global keyframe and keyframe feature set is created, the chunk data (i.e., intra-chunk features, descriptors, correspondences) can be discarded as it is not needed in the second layer pose alignment.

Global Inter-Chunk Pose Optimization Sparse correspondence search and filtering between global keyframes is analogous to that within a chunk, but on the level of all keyframes and their feature sets. If a global keyframe does not find any matches to previously seen keyframes, it is marked as invalid but kept as a candidate, allowing for re-validation when it finds a match to a keyframe observed in the future. The global pose optimization computes the best global alignments $\{\mathcal{T}_i\}$ for the set of all global keyframes, thus aligning all chunks globally. Again, the same energy minimization approach from Sec. 4.3 is applied using both sparse and dense constraints. Intra-chunk alignment runs after each new global keyframe has found correspondences. The pose for a global keyframe is initialized with the delta transform computed by the corresponding intra-chunk optimization, composed with the previous global keyframe pose. After the intra-chunk transforms have been computed, we obtain globally consistent transforms among all input frames by applying the corresponding delta transformations (from the local optimization) to all frames in a chunk.

4.3 Pose Alignment as Energy Optimization

Given a set of 3D correspondences between a set of frames \mathbf{S} (frames in a chunk or keyframes, depending on hierarchy level), the goal of pose alignment is to find an optimal set of rigid camera transforms $\{\mathcal{T}_i\}$ per frame i (for simpler notation, we henceforth write i for f_i) such that all frames align as best as possible. We parameterize the 4×4 rigid transform \mathcal{T}_i using matrix exponentials based on skew-symmetric matrix generators [Murray et al. 1994], which yields fast convergence. This leaves 3 unknown parameters for rotation, and 3 for translation. For ease of notation, we stack the degrees of freedom

for all $|\mathbf{S}|$ frames in a parameter vector:

$$\mathcal{X} = (\mathbf{R}_0, \mathbf{t}_0, \dots, \mathbf{R}_{|\mathbf{S}|}, \mathbf{t}_{|\mathbf{S}|})^T = (x_0, \dots, x_N)^T.$$

Here, N is the total number of variables x_i . Given this notation, we phrase the alignment problem as a variational non-linear least squares minimization problem in the unknown parameters \mathcal{X} . To this end, we define the following alignment objective, which is based on sparse features *and* dense photometric and geometric constraints:

$$E_{\text{align}}(\mathcal{X}) = w_{\text{sparse}} E_{\text{sparse}}(\mathcal{X}) + w_{\text{dense}} E_{\text{dense}}(\mathcal{X}).$$

Here, w_{sparse} and w_{dense} are weights for the sparse and dense matching terms, respectively. w_{dense} is iteratively increased to achieve coarse-to-fine alignment. Note that depending on the optimization hierarchy level, the reference frame is the first frame in the chunk (for intra-chunk alignment), or the first frame in the entire input sequence (for global inter-chunk alignment). Hence, the reference transform \mathcal{T}_0 is not a free variable and left out from the optimization.

Sparse Matching In the sparse matching term, we minimize the sum of distances between the world space positions over all feature correspondences between all pairs of frames in \mathbf{S} :

$$E_{\text{sparse}}(\mathcal{X}) = \sum_{i=1}^{|\mathbf{S}|} \sum_{j=1}^{|\mathbf{S}|} \sum_{(k,l) \in \mathbf{C}(i,j)} \|\mathcal{T}_i \mathbf{p}_{i,k} - \mathcal{T}_j \mathbf{p}_{j,l}\|^2.$$

Here, $\mathbf{p}_{i,k}$ is the k -th detected feature point in the i -th frame. $\mathbf{C}_{i,j}$ is the set of all pairwise correspondences between the i -th and the j -th frame. Geometrically speaking, we seek the best rigid transformations \mathcal{T}_i such that the Euclidean distance over all the detected feature matches is minimized.

Dense Matching We additionally use dense photometric and geometric constraints for fine-scale alignment. To this end, we exploit the dense pixel information of each input frame’s color \mathcal{C}_i and depth \mathcal{D}_i . Evaluating the dense alignment is computationally more expensive than the previous sparse term. We therefore evaluate it on a restricted set \mathbf{E} of frame pairs. \mathbf{E} can be thought of as encoding the edges (i, j) of a sparse matching graph, in which a frame i is aligned (connected with) frame j if their camera angles are similar (within 60° , to avoid glancing angles of the same view) and they overlap with each other. The optimization for both dense photometric and geometric alignment is based on the following energy:

$$E_{\text{dense}}(\mathcal{T}) = w_{\text{photo}} E_{\text{photo}}(\mathcal{T}) + w_{\text{geo}} E_{\text{geo}}(\mathcal{T}).$$

Here, w_{photo} is the weight of the photometric term and w_{geo} of the geometric term, respectively. For the dense photo-consistency term, we evaluate the error on the gradient \mathcal{I}_i of the luminance of \mathcal{C}_i to gain robustness against lighting changes:

$$E_{\text{photo}}(\mathcal{X}) = \sum_{(i,j) \in \mathbf{E}} \sum_{k=0}^{|\mathcal{I}_i|} \left\| \mathcal{I}_i(\pi(\mathbf{d}_{i,k})) - \mathcal{I}_j(\pi(\mathcal{T}_j^{-1} \mathcal{T}_i \mathbf{d}_{i,k})) \right\|_2^2.$$

Here, π denotes the perspective projection, and $\mathbf{d}_{i,k}$ is the 3D position associated with the k -th pixel of the i -th depth frame. Our geometric alignment term evaluates a point-to-plane metric to allow for fine-scale alignment in the tangent plane of the captured geometry:

$$E_{\text{geo}}(\mathcal{X}) = \sum_{(i,j) \in \mathbf{E}} \sum_{k=0}^{|\mathcal{D}_i|} \left[\mathbf{n}_{i,k}^T (\mathbf{d}_{i,k} - \mathcal{T}_i^{-1} \mathcal{T}_j \pi^{-1} (\mathcal{D}_j (\pi(\mathcal{T}_j^{-1} \mathcal{T}_i \mathbf{d}_{i,k})))) \right]^2.$$

Here, $\mathbf{n}_{i,k}$ is the normal of the k -th pixel in the i -th input frame. Correspondences that project outside of the input frame are ignored, and

we apply ICP-like pruning based on distance and normal constraints after each optimization step. For the dense photometric and geometric constraints, we downsample \mathcal{I}_i and \mathcal{D}_i , to 80×60 pixels (using the same cached frames as for the dense verification filter). Note that for the global pose optimization, the result of optimizing densely at every keyframe is effectively reset by the sparse correspondence optimization, since the 3D positions of the correspondences are fixed. Thus we only perform the dense global keyframe optimization after the user has indicated the end of scanning.

4.4 Fast and Robust Optimization Strategy

The described global pose alignment objective is a non-linear least squares problem in the unknown extrinsic camera parameters. Since our goal is *online* and *global* camera pose optimization for long scanning sequences with over twenty thousand frames, an efficient, yet effective, optimization strategy is required. To face this challenge, we implement a data-parallel GPU-based non-linear iterative solver similar to the work of Zollhöfer et al. [2014]. However, the unique sparsity pattern associated with the global alignment objective requires a different parallelization strategy and prohibits the use of previous GPU-based solvers [Zollhöfer et al. 2014; Wu et al. 2014; Zollhöfer et al. 2015]. Our approach is based on the Gauss-Newton method, which only requires first order derivatives and exhibits quadratic convergence close to the optimum, which is beneficial due to our incremental optimization scheme. We find the best pose parameters \mathcal{X}^* by minimizing the proposed highly non-linear least squares objective using this method:

$$\mathcal{X}^* = \underset{\mathcal{X}}{\operatorname{argmin}} E_{align}(\mathcal{X}).$$

For ease of notation, we reformulate the objective in the following canonical least-squares form:

$$E_{align}(\mathcal{X}) = \sum_{i=1}^R r_i(\mathcal{X})^2.$$

This is done by re-naming the $R = 3N_{corr} + |\mathbf{E}| \cdot (|\mathcal{I}_i| + |\mathcal{D}_i|)$ terms of the energy appropriately. Here, N_{corr} is either the total number of inter-chunk sparse correspondences for inter-chunk alignment, or per-chunk sparse correspondences for intra-chunk alignment. The notation can be further simplified by defining a vector field $\mathbf{F} : \mathbb{R}^N \rightarrow \mathbb{R}^R$ that stacks all scalar residuals:

$$\mathbf{F}(\mathcal{X}) = [\dots, r_i(\mathcal{X}), \dots]^T.$$

With this notation, E_{refine} can be expressed in terms of the squared Euclidean length of $\mathbf{F}(\mathcal{X})$:

$$E_{refine}(\mathcal{X}) = \|\mathbf{F}(\mathcal{X})\|_2^2.$$

Gauss-Newton is applied via a local linear approximation of \mathbf{F} at the last solution \mathcal{X}^{k-1} using first-order Taylor expansion:

$$\mathbf{F}(\mathcal{X}^k) = \mathbf{F}(\mathcal{X}^{k-1}) + \mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1}) \cdot \Delta\mathcal{X}, \quad \Delta\mathcal{X} = \mathcal{X}^k - \mathcal{X}^{k-1}.$$

Here, $\mathbf{J}_{\mathbf{F}}$ denotes the Jacobian of \mathbf{F} . By substituting \mathbf{F} with this local approximation, the optimal parameter update $\Delta\mathcal{X}^*$ is found by solving a linear least squares problem:

$$\Delta\mathcal{X}^* = \underset{\Delta\mathcal{X}}{\operatorname{argmin}} \underbrace{\|\mathbf{F}(\mathcal{X}^{k-1}) + \mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1}) \cdot \Delta\mathcal{X}\|_2^2}_{E_{lin}(\Delta\mathcal{X})}.$$

To obtain the minimizer $\Delta\mathcal{X}^*$, we set the corresponding partial derivatives $\frac{dE_{lin}}{d\Delta\mathcal{X}_i}(\Delta\mathcal{X}_i^*) = 0$, $\forall i$ to zero, which yields the following system of linear equations:

$$\mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1})^T \mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1}) \cdot \Delta\mathcal{X}^* = -\mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1})^T \mathbf{F}(\mathcal{X}^{k-1}).$$

To solve the system, we use a GPU-based data-parallel Preconditioned Conjugate Gradient (PCG) solver with Jacobi preconditioner. Based on the iterative solution strategy, the sparsity of the system matrix $\mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1})^T \mathbf{J}_{\mathbf{F}}(\mathcal{X}^{k-1})$ can be exploited. For the sparse term, we never explicitly compute this matrix, but compute the non-zero entries, if required, on-the-fly during the PCG iterations.

Gauss-Newton iterates this process of locally linearizing the energy and solving the associated linear least squares problem starting from an initial estimate \mathcal{X}_0 until convergence. We warm-start the optimization based on the result obtained in the last frame.

In contrast to Zollhöfer [2014], instead of using a reduction based on two kernels to compute the optimal step size and update for the descent direction, we use a single kernel running a combination of warp scans based on the *shuffle* intrinsic and global memory atomics to accumulate the final result. This turned out to be much faster for our problem size.

The central operation of the PCG algorithm is the multiplication of the system matrix with the current descent direction.

Let us first consider the sparse feature term. To avoid fill-in, we multiply the system matrix incrementally based on two separate kernel calls: the first kernel multiplies $\mathbf{J}_{\mathbf{F}}$ and the computed intermediate result is then multiplied by $\mathbf{J}_{\mathbf{F}}^T$ in the second kernel call. For example, at the end of the *apt0* sequence (see Fig. 3, bottom), $\mathbf{J}_{\mathbf{F}}$ has about 105K rows (residuals) and 5K columns (unknowns). In all operations, we exploit the sparse structure of the matrix, only performing operations which will lead to a non-zero result. Since $\mathbf{J}_{\mathbf{F}}$ and $\mathbf{J}_{\mathbf{F}}^T$ have very different row-wise sparsity patterns, using two different kernel calls helps to fine tune the parallelization approach to the specific requirements.

More specifically, for the sparse term, each row of $\mathbf{J}_{\mathbf{F}}$ encodes exactly one pairwise correspondence, depending on at most 2 extrinsic camera poses or $2 \times 6 = 12$ non-zero matrix entries. Due to the low number of required operations, the matrix-vector product can be readily computed by assigning one dedicated thread to each 3D block row; i.e., handling the x -, y -, and z -residuals of one correspondence. This is beneficial, since the different dimensions share common operations in the evaluation of \mathbf{F} and $\mathbf{J}_{\mathbf{F}}$. In contrast, \mathbf{J}^T has exactly one row per unknown. The number of non-zero entries in each row is equivalent to the number of correspondences involving the frame associated with the unknown. For longer scanning sequences, this can easily lead to several thousand entries per row. To reduce the amount of memory reads and compute of each thread, we opted for a reduction-based approach to compute the matrix-vector products. We use one block of size $N_{block} = 256$ to compute each row-wise dot product. Each warp of a block performs a warp-reduction based on the *shuffle* intrinsic and the final per-warp results are combined based on shared memory atomics. For computing the multiplication with $\mathbf{J}_{\mathbf{F}}^T$, we pre-compute auxiliary lists that allow lookup to all correspondences that influence a certain variable. This table is filled based on a kernel that has one thread per correspondence and adds entries to the lists corresponding to the involved variables. The per-list memory is managed using *atomic* counters. We recompute this table if the set of active correspondences changes.

For the dense photometric and geometric alignment terms, the number of associated residuals is considerably higher. Since the system matrix is fixed during the PCG steps, we pre-compute it at the beginning of each non-linear iteration. The required memory is preallocated and we update only the non-zero entries via scattered writes. Note that we only require a few writes, since we perform the local reductions in shared memory.

4.4.1 Correspondence and Frame Filtering. As an additional safeguard to make the optimization robust against potential corre-

spendence outliers, which were mistakenly considered to be valid, we perform correspondence and frame filtering after each optimization finishes. That is, we determine the maximum residual $r_{max} = \max_i r_i(\mathcal{X})$ using a parallel reduction on the GPU, with the final max computation performed on the CPU. If $r_{max} > 0.05m$, we remove all correspondences between the two frames i and j associated with the correspondence which induces r_{max} . Note that all correspondences between i and j are removed, in order to minimize the number of times the optimization has to run in order to prune all bad correspondences. Additionally, if a frame has no correspondence to any other frame, it is implicitly removed from the optimization and marked as invalid.

5. DYNAMIC 3D RECONSTRUCTION

Key to live, globally consistent reconstruction is updating the 3D model based on newly-optimized camera poses. We thus monitor the continuous change in the poses of each frame to update the volumetric scene representation through *integration* and *de-integration* of frames. Based on this strategy, errors in the volumetric representation due to accumulated drift or dead reckoning in feature-less regions can be fixed as soon as better pose estimates are available.

5.1 Scene Representation

Scene geometry is reconstructed by incrementally fusing all input RGB-D data into an implicit truncated signed distance (TSDF) representation, following Curless and Levoy [1996]. The TSDF is defined over a volumetric grid of voxels; to store and process this data, we employ the state-of-the-art sparse volumetric voxel hashing approach proposed by Nießner et al. [2013]. This approach scales well to the scenario of large-scale surface reconstruction, since empty space neither needs to be represented nor addressed; the TSDF is stored in a sparse volumetric grid based on spatial hashing. Following the original approach, we also use voxel blocks of $8 \times 8 \times 8$ voxels. In contrast to the work of Nießner et al. [2013], we allow for RGB-D frames to both be *integrated* into the TSDF as well as *de-integrated* (i.e., adding and removing frames from the reconstruction). In order to allow for pose updates, we also ensure that these two operations are symmetric; i.e., one inverts the other.

5.2 Integration and De-integration

Integration of a depth frame \mathcal{D}_i occurs as follows. For each voxel, $\mathbf{D}(\mathbf{v})$ denotes the signed distance of the voxel, $\mathbf{W}(\mathbf{v})$ the voxel weight, $d_i(\mathbf{v})$ the projective distance (along the z axis) between a voxel and \mathcal{D}_i , and $w_i(\mathbf{v})$ the integration weight for a sample of \mathcal{D}_i . For data integration, each voxel is then updated by

$$\mathbf{D}'(\mathbf{v}) = \frac{\mathbf{D}(\mathbf{v})\mathbf{W}(\mathbf{v}) + w_i(\mathbf{v})d_i(\mathbf{v})}{\mathbf{W}(\mathbf{v}) + w_i(\mathbf{v})}, \quad \mathbf{W}'(\mathbf{v}) = \mathbf{W}(\mathbf{v}) + w_i(\mathbf{v}).$$

We can reverse this operation to de-integrate a frame. Each voxel is then updated by

$$\mathbf{D}'(\mathbf{v}) = \frac{\mathbf{D}(\mathbf{v})\mathbf{W}(\mathbf{v}) - w_i(\mathbf{v})d_i(\mathbf{v})}{\mathbf{W}(\mathbf{v}) - w_i(\mathbf{v})}, \quad \mathbf{W}'(\mathbf{v}) = \mathbf{W}(\mathbf{v}) - w_i(\mathbf{v}).$$

We can thus update a frame in the reconstruction by de-integrating it from its original pose and integrating it with a new pose. This is crucial for obtaining high-quality reconstructions in the presence of loop closures and revisiting, since the already integrated surface measurements must be adapted to the continuously changing stream of pose estimates.

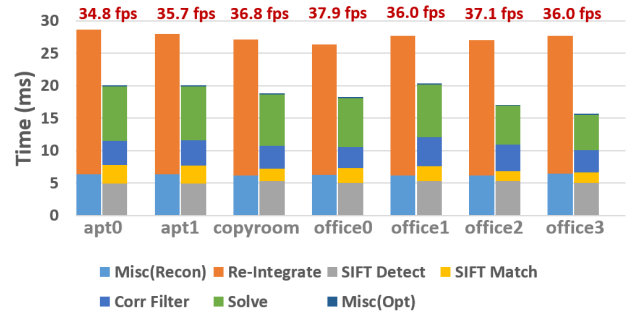


Fig. 4. Performance Evaluation: our proposed pipeline runs at well beyond 30Hz for all used test sequences. The computations are split up over two GPUs (left bar Titan X, right bar Titan Black).

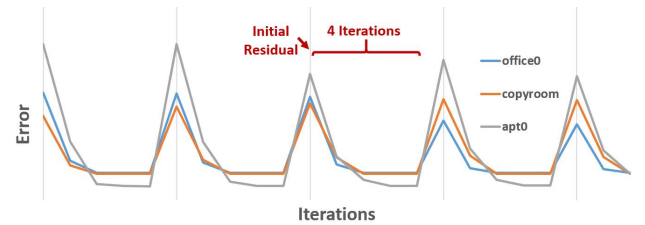


Fig. 5. Convergence analysis of the global keyframe optimization (log scale): peaks correspond to new global keyframes. Only a few iterations are required for convergence.

5.3 Managing Reconstruction Updates

Each input frame is stored with its associated depth and color data, along with two poses: its *integrated* pose, and its *optimized* pose. The integrated pose is the one used currently in the reconstruction, and is set whenever a frame gets integrated. The optimized pose stores the (continually-changing) result of the pose optimization.

When an input frame arrives, we aim to integrate it into the reconstruction as quickly as possible, to give the user or robot instantaneous feedback of the 3D model. Since the global optimization is not run for each frame but for each chunk, an optimized pose may not immediately be available and we must obtain an initial transform by other means. We compute this initial transform by composing the frame-to-frame transform from the key point correspondence filter with the newest available optimized transform.

In order to update the reconstruction with the most pertinent optimization updates, we sort the frames in descending order by the difference between the integrated transform and the optimized transform. The integrated transform and optimized transform are parameterized by 6 DOFs: α, β, γ (here, we use Euler angles in radians) describing the rotation, and x, y, z (in meters) describing the translation. Then the distance between the integrated transform $\mathbf{t}_{int} = (\alpha_i, \beta_i, \gamma_i, x_i, y_i, z_i)$ and the optimized transform $\mathbf{t}_{opt} = (\alpha_o, \beta_o, \gamma_o, x_o, y_o, z_o)$ is defined to be $\|\mathbf{s} * \mathbf{t}_{int} - \mathbf{s} * \mathbf{t}_{opt}\|_2$ where $\mathbf{s} = (2, 2, 2, 1, 1, 1)$ is multiplied element-wise to bring the rotations and translations closer in scale. For each new input frame, we de-integrate and integrate the $N_{fix} = 10$ frames from the top of the list. This allows us to dynamically update the reconstruction to produce a globally-consistent 3D reconstruction.



Fig. 3. Large-scale reconstruction results: our proposed real-time global pose optimization outperforms current state-of-the-art online reconstruction systems. The globally aligned 3D reconstructions are at a quality that was previously only attainable offline. Note the completeness of the scans, the global alignment without noticeable camera drift and the high local quality of the reconstructions in both geometry and texture. Scans comprise thousands of input frames, include revisiting and many loop closures.

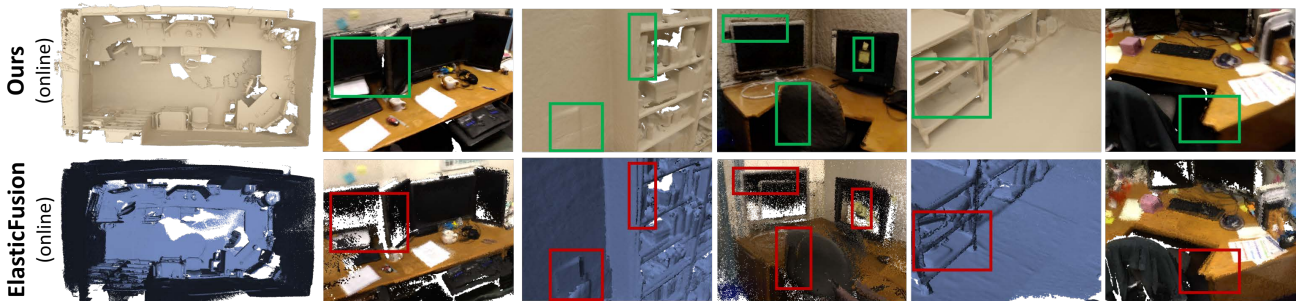


Fig. 7. Our proposed real-time global pose optimization (top) outperforms the method of Whelan et al. [2015] (bottom) in terms of scan completeness and alignment accuracy. Note, we generate a high-quality surface mesh, while the competing approach only outputs a pointcloud.

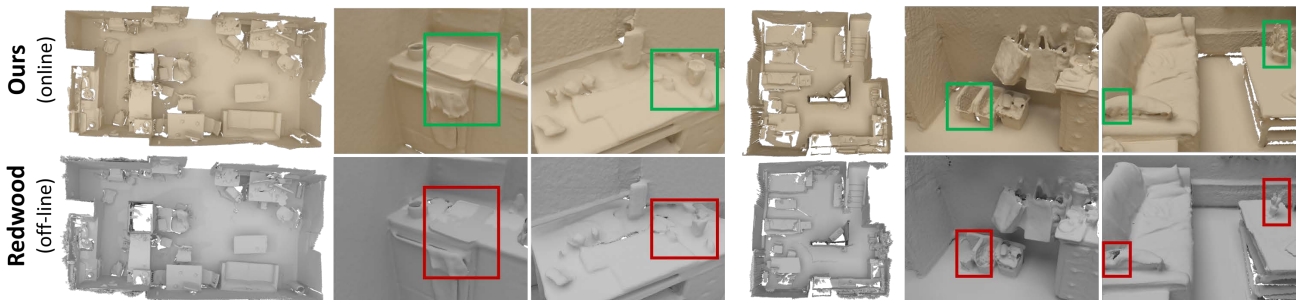


Fig. 8. Our proposed real-time global pose optimization (top) delivers a reconstruction quality on par or even better than the off-line Redwood [Choi et al. 2015] system (bottom). Note, our reconstructions have more small scale detail.

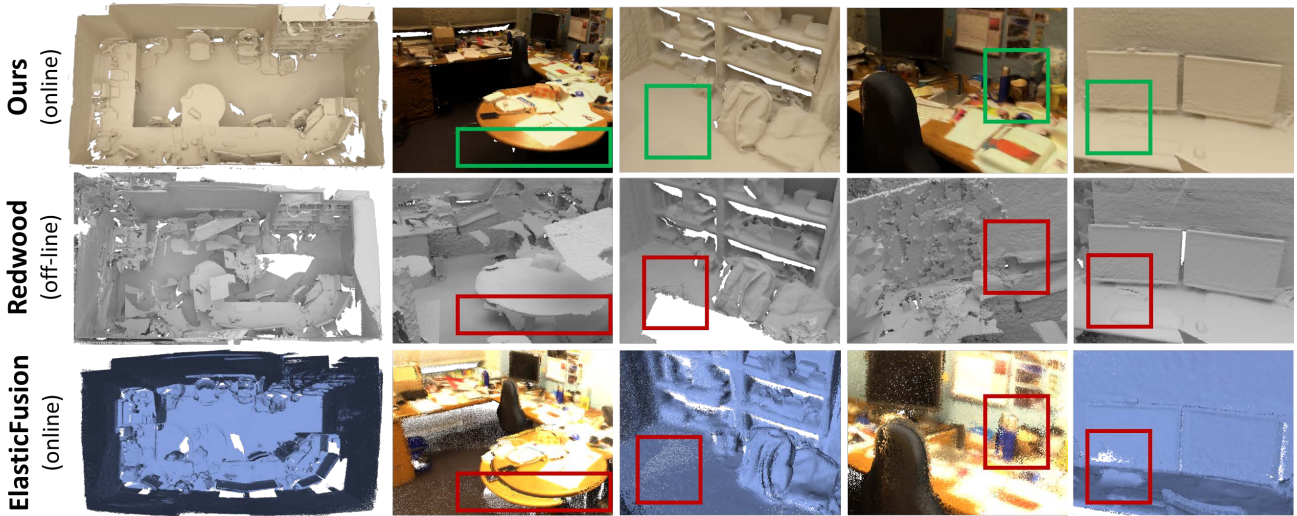


Fig. 9. Our proposed real-time global pose optimization (top) delivers a reconstruction quality on par or even better than the off-line Redwood [Choi et al. 2015] (middle) and the ElasticFusion [2015] (bottom) system. Note that Redwood does not use color information, and was not able to resolve all loop closures in this challenging scan.

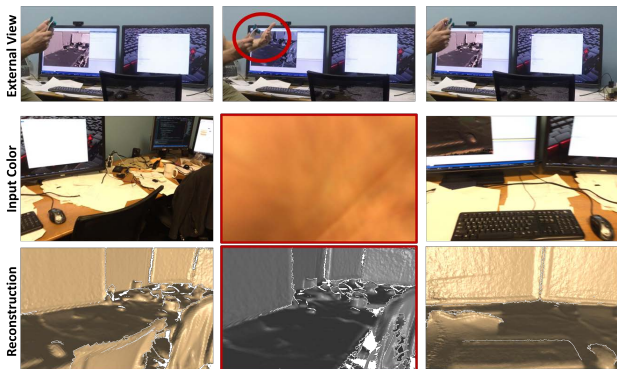


Fig. 6. Recovery from tracking failure: our method is able to detect (gray overlay) and recover from tracking failure; i.e., if the sensor is occluded or observes a featureless region.

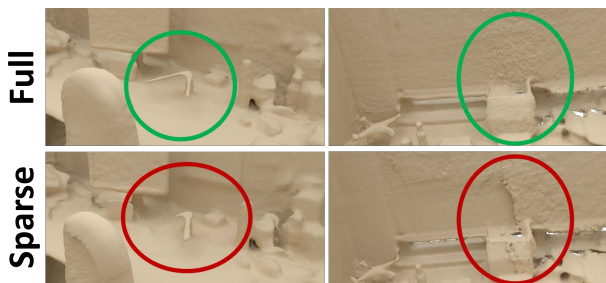


Fig. 10. Comparison of Sparse vs. Dense Alignment: the proposed dense intra- and inter- chunk alignment (top) leads to higher quality reconstructions than only the sparse alignment step (bottom).

6. RESULTS

For live scanning, we use a *Structure Sensor*² mounted to an iPad Air. The RGB-D stream is captured at 30Hz with a color and depth resolution of 640×480 . Note that we are agnostic to the type of used depth sensor. We stream the captured RGB-D data via a wireless network connection to a desktop machine that runs our global pose optimization and reconstructs a 3D model in real-time. Visual feedback of the reconstruction is streamed live to the iPad to aid in the scanning process. To reduce the required bandwidth, we use data compression based on *zlib* for depth and *jpeg* compression for color. We implemented our global pose alignment framework using the CUDA 7.0 architecture. Reconstruction results of scenes captured using our live system are shown in Fig. 1 and 3 as well as in the supplementary video. The completeness of the various large-scale indoor scenes (4 offices, 2 apartments, 1 copyroom, with up to 95m camera trajectories), their alignment without noticeable camera drift, and the high local quality of geometry and texture are on par with even *offline* approaches. This also demonstrates that our global pose alignment strategy scales well to large spatial extents and long sequences (over 20,000 frames).

Qualitative Comparison. First, we compare to the online 3D reconstruction approach of Nießner et al. [2013], see Fig. 11. In contrast to their work, which builds on frame-to-model tracking and suffers from the accumulation of camera drift, we are able to produce drift-free reconstructions at high fidelity. Our novel global pose optimization framework implicitly handles loop closure, recovers from tracking failures, and reduces geometric drift. Note that most real-time fusion methods (e.g., [Izadi et al. 2011; Newcombe et al. 2011; Chen et al. 2013; Nießner et al. 2013]) share the same frame-to-model ICP tracking algorithm, and therefore suffer from notable drift. Fig. 7 and 9 show a comparison of our approach with the online *ElasticFusion* approach of Whelan et al. [2015], which captures surfel maps using dense frame-to-model tracking and explicitly

²<http://structure.io/>

handles loop closures using non-rigid warping. In contrast, our dynamic de-integration and integration of frames mitigates issues with warping artifacts in rigid structures, and moreover produces a high quality continuous surface. Since our approach does not rely on explicit loop closure detection, it scales better to scenarios with many loop closures (c.p. Fig. 7 and 9). We additionally compare to the offline Redwood approach [Choi et al. 2015], using their rigid variant, see Fig. 8 and 9. Note, we do not compare to their newer non-rigid approach, since it fails on most of our dataset sequences. While their approach takes several hours (2.3h - 13.2h for each of our sequences), we achieve comparable quality and better reconstruction of small-scale detail at real-time rates. Note that Redwood does not take color information into account, thus struggling with sequences that contain fewer geometric features.

Performance and Convergence. We measure the performance of our pipeline on an Intel Core i7 3.4GHz CPU (32GB RAM). For compute, we use a combination of a NVIDIA GeForce GTX Titan X and a GTX Titan Black. The Titan X is used for volumetric reconstruction, and the Titan Black for correspondence search and global pose optimization. Our pipeline runs with a framerate well beyond 30Hz (see Fig. 4) for all shown test sequences. Note that the global dense optimization runs in < 500 ms at the end of the sequences. After adding a new global keyframe, our approach requires only a few iterations to reach convergence. Fig. 5 shows convergence plots for three of the used test sequences (cf. Fig. 3); the behavior generalizes to all other sequences.

Recovery from Tracking Failure. If a new keyframe cannot be aligned successfully, we assume tracking is lost and do not integrate surface measurements. An example scanning sequence is shown in Fig. 6. To indicate tracking failure, the reconstruction is shown with a gray overlay. Based on this cue, the user is able to recover the method by moving back to a previously scanned area. Note that there is no temporal nor spatial coherence required, as our method globally matches new frames against all existing data. Thus, scanning may be interrupted, and continued at a completely different location.

Loop Closure Detection and Handling. Our global pose optimization approach detects and handles loop closures transparently (see Fig. 12), since the volumetric scene representation is continuously updated to match the stream of computed pose estimates. This allows incrementally fixing of loop closures over time by means of *integration* and *de-integration* of surface measurements.

Dense Tracking and Voxel Resolution. In Fig. 10, we evaluate the influence of the dense tracking component of our energy function. While globally drift-free reconstructions can be obtained by sparse tracking only, the dense alignment term leads to more refined local results. The impact of voxel resolution on reconstruction quality is shown in Fig. 13. As a default, we use a voxel resolution of 4mm for all reconstructions. While 1cm voxels reduce memory consumption, the quality of the reconstruction is slightly impaired.

Quantitative Comparison. We quantitatively evaluate our approach on independent benchmark data and compare against state-of-the-art online (DVO-SLAM [Kerl et al. 2013], RGB-D SLAM [Endres et al. 2012], MRSMap [Stückler and Behnke 2014], Kintinuous [Whelan et al. 2012], VoxelHashing [Nießner et al. 2013; Nießner et al. 2014], ElasticFusion [Whelan et al. 2015]) and offline systems (Submap Bundle Adjustment [Maier et al. 2014], Redwood [Choi et al. 2015]). Note that for Redwood, we show results for the rigid variant, which produced better camera tracking results. We first evaluate our approach on the ICL-NUIM dataset of Handa

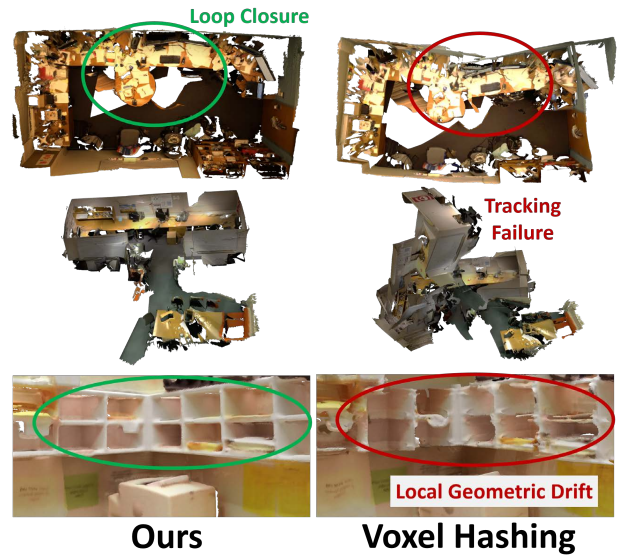


Fig. 11. Comparison to the VoxelHashing approach of Nießner et al. [2013]: in contrast to the frame-to-model tracking of VoxelHashing, our novel global pose optimization implicitly handles loop closure (top), robustly detects, recovers from tracking failures (middle), and greatly reduces geometric drift (bottom).

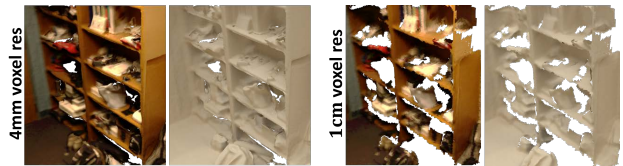


Fig. 13. Comparison of different voxel resolutions: 4mm voxel resolution (left) leads to higher-fidelity reconstructions than the coarser 1cm resolution (right). Note the generally sharper texture and the more refined geometry in case of 4mm voxels.

et al. [2014], which provides ground truth camera poses for several scans of a synthetic environment. Table I shows our trajectory estimation performance, measured with absolute trajectory error (ATE), on the four living room scenes (including synthetic noise), for which we out-perform existing state-of-the-art online and offline systems. Additionally, in Table II we evaluate our approach on the RGB-D benchmark of Sturm et al [2012]. This benchmark provides ground truth camera pose estimates for hand-held Kinect sequences using a calibrated motion capture system. For these sequences, which only cover small scenes and simple camera trajectories, our results are on par with or better than the existing state of the art. Note that our own sequences have a larger spatial extent and are much more challenging, with faster motion and many more loop closures. For these data sets (Tables I and II), the Redwood system, which relies solely on geometric registration, suffers from the relative lack of varying views in the camera trajectories. In particular, fr3/nst is a textured wall, which cannot be registered with a geometric-only method. On both these datasets, we also quantitatively validate the relevance of our design decisions. While online alignment based on sparse features only (Ours (s)) achieves reasonable results, using dense matching only in per chunk alignment further increases accuracy

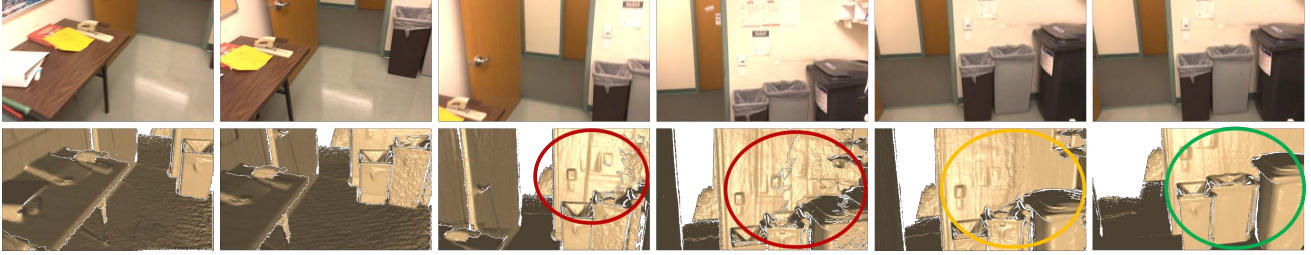


Fig. 12. Detection and resolving of loop closure is enabled through global pose optimization. Note, while data is first integrated at slightly wrong locations, the volumetric representation is fixed over time as soon as better pose estimates are available.

(Ours (sd)). Our full sparse and dense matching approach on both local and global level leads to highest accuracy.

Table I. ATE RMSE on the synthetic ICL-NUIM dataset by [Handa et al. 2014].

	kt0	kt1	kt2	kt3
DVO SLAM	10.4cm	2.9cm	19.1cm	15.2cm
RGB-D SLAM	2.6cm	0.8cm	1.8cm	43.3cm
MRSMap	20.4cm	22.8cm	18.9cm	109cm
Kintinuous	7.2cm	0.5cm	1.0cm	35.5cm
VoxelHashing	1.4cm	0.4cm	1.8cm	12.0cm
Elastic Fusion	0.9cm	0.9cm	1.4cm	10.6cm
Redwood (rigid)	25.6cm	3.0cm	3.3cm	6.1cm
Ours (s)	0.9cm	1.2cm	1.3cm	1.3cm
Ours (sd)	0.8cm	0.5cm	1.1cm	1.2cm
Ours	0.6cm	0.4cm	0.6cm	1.1cm

Note that unlike the other methods, Redwood does not use color information and runs offline. For our approach, we also provide results for sparse-only (s) as well as sparse and local dense only (sd).

Table II. ATE RMSE on the TUM RGB-D dataset by [Sturm et al. 2012].

	fr1/desk	fr2/xyz	fr3/office	fr3/nst
DVO SLAM	2.1cm	1.8cm	3.5cm	1.8cm
RGB-D SLAM	2.3cm	0.8cm	3.2cm	1.7cm
MRSMap	4.3cm	2.0cm	4.2cm	201.8cm
Kintinuous	3.7cm	2.9cm	3.0cm	3.1cm
VoxelHashing	2.3cm	2.2cm	2.3cm	8.7cm
Elastic Fusion	2.0cm	1.1cm	1.7cm	1.6cm
LSD-SLAM	-	1.5cm	-	-
Submap BA	2.2cm	-	3.5cm	-
Redwood (rigid)	2.7cm	9.1cm	3.0cm	192.9cm
Ours (s)	1.9cm	1.4cm	2.9cm	1.6cm
Ours (sd)	1.7cm	1.4cm	2.8cm	1.4cm
Ours	1.6cm	1.1cm	2.2cm	1.2cm

Note that unlike the other methods listed, Redwood does not use color information and runs offline. For our approach, we also provide results for sparse-only (s) as well as sparse and local dense only (sd).

Limitations. As our tracking is based on sparse key point matching, small local misalignments can occur; e.g., SIFT matches can be off by a few pixels and the depth data associated with a keypoint may be inaccurate due to sensor noise. While we solve for optimal keypoint positions from the inter-chunk optimization, small mismatches between global keypoints can still be propagated within the global

optimization, leading to local misalignments. Ideally, we would treat the locations of the global keypoints as unknowns to optimize for. Unfortunately, this would involve significant computational effort, which (currently) seems to exceed even the computational budget of offline approaches. Another limitation is that we currently run our method on two GPUs. Fortunately, we can easily stream the data to and from an iPad with live visual feedback, on both the desktop and mobile device, thus making scanning fun and convenient.

7. CONCLUSION

We have presented a novel online real-time 3D reconstruction approach that provides robust tracking and implicitly solves the loop closure problem by globally optimizing the trajectory for every captured frame. To this end, we combine online SIFT feature extraction, matching, and pruning with a novel parallel non-linear pose optimization framework, over both sparse features as well as dense correspondences, enabling the solution of the global alignment problem at real-time rates. The continuously changing stream of optimized pose estimates is monitored and the reconstruction is updated through dynamic integration and de-integration. The capabilities of the proposed approach have been demonstrated on several large-scale 3D reconstructions with a reconstruction quality and completeness that was previously only possible with offline approaches and tedious capture sessions. We believe online global pose alignment will pave the way for many new and interesting applications. Global accurate tracking is the foundation for immersive AR/VR applications and makes online hand-held 3D reconstruction applicable to scenarios that require high-fidelity tracking.

ACKNOWLEDGMENTS

We would like to thank Thomas Whelan for his help with Elastic-Fusion, and Sungjoon Choi for his advice on the Redwood system. This work was funded by the Max Planck Center for Visual Computing and Communications, the ERC Starting Grant 335545 CapReal, and a Stanford Graduate Fellowship. We are grateful for hardware donations from NVIDIA Corporation.

REFERENCES

- BESL, P. AND MCKAY, N. 1992. A method for registration of 3-D shapes. *IEEE Trans. PAMI* 14, 2, 239–256.
- CHEN, J., BAUTEMBACH, D., AND IZADI, S. 2013. Scalable real-time volumetric surface reconstruction. *ACM TOG* 32, 4, 113.
- CHOI, S., ZHOU, Q.-Y., AND KOLTUN, V. 2015. Robust reconstruction of indoor scenes. *Proc. CVPR*.

- CURLLESS, B. AND LEVOY, M. 1996. A volumetric method for building complex models from range images. In *In Proc. SIGGRAPH*. ACM, 303–312.
- ELFES, A. AND MATTHIES, L. 1987. Sensor integration for robot navigation: combining sonar and stereo range data in a grid-based representation. In *Decision and Control, 1987. 26th IEEE Conference on*. Vol. 26. IEEE, 1802–1807.
- ENDRES, F., HESS, J., ENGELHARD, N., STURM, J., CREMERS, D., AND BURGARD, W. 2012. An evaluation of the rgb-d slam system. In *Proc. ICRA*. IEEE, 1691–1696.
- ENGEL, J., SCHÖPS, T., AND CREMERS, D. 2014. LSD-SLAM: Large-scale direct monocular SLAM.
- ENGEL, J., STURM, J., AND CREMERS, D. 2013. Semi-dense visual odometry for a monocular camera. In *Proc. ICCV*. IEEE, 1449–1456.
- FIORAI, N., TAYLOR, J., FITZGIBBON, A., DI STEFANO, L., AND IZADI, S. 2015. Large-scale and drift-free surface reconstruction using online subvolume registration. *Proc. CVPR*.
- FORSTER, C., PIZZOLI, M., AND SCARAMUZZA, D. 2014. Svo: Fast semi-direct monocular visual odometry. In *Proc. ICRA*. IEEE, 15–22.
- FUHRMANN, S. AND GOESELE, M. 2014. Floating Scale Surface Reconstruction. In *Proc. SIGGRAPH*.
- GALLUP, D., POLLEFEYS, M., AND FRAHM, J.-M. 2010. 3D reconstruction using an n-layer heightmap. In *Pattern Recognition*. Springer, 1–10.
- GLOCKER, B., SHOTTON, J., CRIMINISI, A., AND IZADI, S. 2015. Real-time rgb-d camera relocalization via randomized ferns for keyframe encoding. *TVCG 21*, 5, 571–583.
- GOWER, J. C. 1975. Generalized procrustes analysis. *Psychometrika* 40, 1, 33–51.
- HANDA, A., WHELAN, T., McDONALD, J., AND DAVISON, A. 2014. A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM. In *Proc. ICRA*. Hong Kong, China.
- HENRY, P., KRAININ, M., HERBST, E., REN, X., AND FOX, D. 2010. RGB-D mapping: Using depth cameras for dense 3d modeling of indoor environments. In *Proc. Int. Symp. Experimental Robotics*. Vol. 20. 22–25.
- HILTON, A., STODDART, A., ILLINGWORTH, J., AND WINDEATT, T. 1996. Reliable surface reconstruction from multiple range images. *JProc. ECCV*, 117–126.
- IZADI, S., KIM, D., HILLIGES, O., MOLYNEAUX, D., NEWCOMBE, R., KOHLI, P., SHOTTON, J., HODGES, S., FREEMAN, D., DAVISON, A., AND FITZGIBBON, A. 2011. KinectFusion: Real-time 3D reconstruction and interaction using a moving depth camera. In *Proc. UIST*. 559–568.
- KABSCH, W. 1976. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 32, 5, 922–923.
- 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 *Proc. 3DV*. IEEE, 1–8.
- KERL, C., STURM, J., AND CREMERS, D. 2013. Dense visual slam for rgb-d cameras. In *Proc. IROS*.
- KLEIN, G. AND MURRAY, D. 2007. Parallel tracking and mapping for small AR workspaces. In *Proc. ISMAR*. Nara, Japan.
- KÜMMERLE, R., GRISETTI, G., STRASDAT, H., KONOLIGE, K., AND BURGARD, W. 2011. g 2 o: A general framework for graph optimization. In *Proc. ICRA*. IEEE, 3607–3613.
- LEVOY, M., PULLI, K., CURLLESS, B., RUSINKIEWICZ, S., KOLLER, D., PEREIRA, L., GINZTON, M., ANDERSON, S., DAVIS, J., GINSBERG, J., ET AL. 2000. The digital michelangelo project: 3D scanning of large statues. In *In Proc. SIGGRAPH*. ACM Press/Addison-Wesley Publishing Co., 131–144.
- LI, H., VOUGA, E., GUDYM, A., LUO, L., BARRON, J. T., AND GUSEV, G. 2013. 3d self-portraits. *ACM TOG* 32, 6, 187.
- LOWE, D. G. 2004. Distinctive image features from scale-invariant keypoints. *IJCV* 60, 91–110.
- MAIER, R., STURM, J., AND CREMERS, D. 2014. Submap-based bundle adjustment for 3d reconstruction from rgb-d data. In *Proc. GCPR*. Münster, Germany.
- MEILLAND, M., COMPORT, A., ET AL. 2013. On unifying key-frame and voxel-based dense visual slam at large scales. In *Proc. IROS*. IEEE, 3677–3683.
- MEILLAND, M., COMPORT, A., RIVES, P., AND MÉDITERRANÉE, I. S. A. 2011. Real-time dense visual tracking under large lighting variations. In *Proc. BMVC*. Vol. 29.
- MERRELL, P., AKBARZADEH, A., WANG, L., MORDOHAJ, P., FRAHM, J., YANG, R., NISTÉR, D., AND POLLEFEYS, M. 2007. Real-time visibility-based fusion of depth maps. In *Proc. ICCV*. 1–8.
- MURRAY, R. M., SASTRY, S. S., AND ZEXIANG, L. 1994. *A Mathematical Introduction to Robotic Manipulation*. CRC Press.
- 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. ISMAR*. 127–136.
- NEWCOMBE, R. A., LOVEGROVE, S. J., AND DAVISON, A. J. 2011. Dtam: Dense tracking and mapping in real-time. In *Proc. ICCV*. 2320–2327.
- NISSNER, M., DAI, A., AND FISHER, M. 2014. Combining inertial navigation and icp for real-time 3d surface reconstruction.
- NISSNER, M., ZOLLHÖFER, M., IZADI, S., AND STAMMINGER, M. 2013. Real-time 3d reconstruction at scale using voxel hashing. *ACM TOG*.
- PRADEEP, V., RHEMANN, C., IZADI, S., ZACH, C., BLEYER, M., AND BATHICHE, S. 2013. Monofusion: Real-time 3d reconstruction of small scenes with a single web camera. In *Proc. ISMAR*. 83–88.
- REICHL, F., WEISS, J., AND WESTERMANN, R. 2015. Memory-efficient interactive online reconstruction from depth image streams. In *Computer Graphics Forum*. Wiley Online Library.
- ROTH, H. AND VONA, M. 2012. Moving volume KinectFusion. In *Proc. BMVC*.
- RUSINKIEWICZ, S., HALL-HOLT, O., AND LEVOY, M. 2002. Real-time 3D model acquisition. *ACM TOG* 21, 3, 438–446.
- RUSINKIEWICZ, S. AND LEVOY, M. 2001. Efficient variants of the ICP algorithm. In *Proc. 3DIM*. 145–152.
- SILBERMAN, N., HOIEM, D., KOHLI, P., AND FERGUS, R. 2012. Indoor segmentation and support inference from rgb-d images. In *ECCV*.
- STEINBRUECKER, F., KERL, C., STURM, J., AND CREMERS, D. 2013. Large-scale multi-resolution surface reconstruction from rgb-d sequences. In *Proc. ICCV*. Sydney, Australia.
- STEINBRUECKER, F., STURM, J., AND CREMERS, D. 2014. Volumetric 3d mapping in real-time on a cpu. Hongkong, China.
- STÜCKLER, J. AND BEHNKE, S. 2014. Multi-resolution surfel maps for efficient dense 3d modeling and tracking. *Journal of Visual Communication and Image Representation* 25, 1, 137–147.
- STURM, J., ENGELHARD, N., ENDRES, F., BURGARD, W., AND CREMERS, D. 2012. A benchmark for the evaluation of rgb-d slam systems. In *Proc. IROS*.
- TRIGGS, B., MCLAUCHLAN, P. F., HARTLEY, R. I., AND FITZGIBBON, A. W. 2000. Bundle adjustment, a modern synthesis. In *Vision algorithms: theory and practice*. Springer, 298–372.
- VALENTIN, J., NISSNER, M., SHOTTON, J., FITZGIBBON, A., IZADI, S., AND TORR, P. 2015. Exploiting uncertainty in regression forests for accurate camera relocalization. In *Proc. CVPR*. 4400–4408.
- WEISE, T., WISMER, T., LEIBE, B., AND VAN GOOL, L. 2009. In-hand scanning with online loop closure. In *Proc. ICCV Workshops*. 1630–1637.

- WHELAN, T., JOHANSSON, H., KAESS, M., LEONARD, J., AND MCDONALD, J. 2012. Robust tracking for real-time dense rgb-d mapping with kintinuous. Tech. rep. Query date: 2012-10-25.
- WHELAN, T., JOHANSSON, H., KAESS, M., LEONARD, J. J., AND MCDONALD, J. 2013. Robust real-time visual odometry for dense rgb-d mapping. In *Proc. ICRA*.
- WHELAN, T., KAESS, M., LEONARD, J. J., AND MCDONALD, J. 2013. Deformation-based loop closure for large scale dense rgb-d slam. In *Proc. IROS*. IEEE, 548–555.
- WHELAN, T., LEUTENEGGER, S., SALAS-MORENO, R. F., GLOCKER, B., AND DAVISON, A. J. 2015. ElasticFusion: Dense SLAM without a pose graph. In *Proc. RSS*. Rome, Italy.
- WU, C., ZOLLHÖFER, M., NIESSNER, M., STAMMINGER, M., IZADI, S., AND THEOBALT, C. 2014. Real-time shading-based refinement for consumer depth cameras. *ACM TOG* 33, 6.
- WURM, K. M., HORNUNG, A., BENNEWITZ, M., STACHNISS, C., AND BURGARD, W. 2010. Octomap: A probabilistic, flexible, and compact 3d map representation for robotic systems. In *Proc. ICRA*. Vol. 2.
- XIAO, J., OWENS, A., AND TORRALBA, A. 2013. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *Proc. ICCV*. IEEE, 1625–1632.
- ZENG, M., ZHAO, F., ZHENG, J., AND LIU, X. 2012. Octree-based fusion for realtime 3D reconstruction. *Graphical Models*.
- ZHOU, Q.-Y. AND KOLTUN, V. 2013. Dense scene reconstruction with points of interest. *ACM Transactions on Graphics (TOG)* 32, 4, 112.
- ZHOU, Q.-Y. AND KOLTUN, V. 2014. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Transactions on Graphics (TOG)* 33, 4, 155.
- ZHOU, Q.-Y., MILLER, S., AND KOLTUN, V. 2013. Elastic fragments for dense scene reconstruction. In *Computer Vision (ICCV), 2013 IEEE International Conference on*. IEEE, 473–480.
- ZOLLHÖFER, M., DAI, A., INNEMANN, M., WU, C., STAMMINGER, M., THEOBALT, C., AND NIESSNER, M. 2015. Shading-based refinement on volumetric signed distance functions. *ACM TOG* 34, 4.
- ZOLLHÖFER, M., NIESSNER, M., IZADI, S., REHMANN, C., ZACH, C., FISHER, M., WU, C., FITZGIBBON, A., LOOP, C., THEOBALT, C., AND STAMMINGER, M. 2014. Real-time non-rigid reconstruction using an rgb-d camera. *ACM TOG* 33, 4.

APPENDIX

Table III. Dataset overview.

	#Frames	Trajectory Length
Apt 0	8560	89.4m
Apt 1	8580	91.8m
Apt 2	3989	87.5m
Copyroom	4480	24.4m
Office 0	6159	52.8m
Office 1	5730	51.7m
Office 2	3500	36.3m
Office 3	3820	66.8m

In order to capture scans at high completeness, the camera is moved in long and complex trajectories.

A. ADDITIONAL QUALITATIVE RESULTS

Reconstructed models for the seven scenes in our dataset are publicly available at <http://www.graphics.stanford.edu/projects/bundlefusion/>. Note that the relocalization (due to sensor occlusion) in the sequence Apt 2 cannot be handled by state-of-the-art methods such as ElasticFusion and Redwood. While our method and ElasticFusion run at real-time rates, Redwood runs offline, taking 8.6 hours for Apt0, 13.2 hours for Apt1, 4 hours for Copyroom, 7.7 hours for Office1, 2.6 hours for Office2, and 3.5 hours for Office3. Redwood is also a geometry-only approach that does not use the RGB channels.

We additionally evaluate our method on the SUN3D dataset [Xiao et al. 2013], which contains a variety of indoor scenes captured with an Asus Xtion sensor. Fig. 14 shows reconstruction results for several large, complex scenes, using the offline SUN3Dsfm bundle adjustment system as well as our approach. Note that our approach produces better global structure while maintaining local detail at real-time rates. The SUN3D dataset also contains eight scenes which contain manual object-correspondence annotations in order to guide their reconstructions; we show reconstruction results using our method (without annotation information) on these scenes in Fig. 15.

We have also reconstructed all 464 scenes from the NYU2 dataset [Silberman et al. 2012], which contains a variety of indoor scenes recorded by a Kinect. Several reconstruction results are shown in Fig. 16.

B. ADDITIONAL QUANTITATIVE RESULTS

The ICL-NUIM dataset of Handa et al. [2014] also provides the ground truth 3D model used to generate the virtually scanned sequences. In addition to the camera tracking evaluation provided in Section 6 of the paper, we evaluate surface reconstruction accuracy (mean distance of the model to the ground truth surface) for the living room model in Table IV.

Additionally, we further evaluate our camera tracking on the augmented ICL-NUIM dataset of [Choi et al. 2015], which comprises synthetic scans of two virtual scenes, a living room and an office, from the original ICL-NUIM data. In contrast to the original ICL-NUIM, these scans have longer trajectories with more loop closures. Table V shows our trajectory estimation performance on this dataset (with synthetic sensor noise, using the reported camera intrinsic parameters), which is on par with or better than existing state of the art. Although the camera trajectories are complex, the additional

loop closures help maintain stability (as frames find matches which are not neighbors, mitigating tracking drift), aiding our performance in all scenes except Office 1. In this case, our method has difficulty closing the loop, as part of it covers a wall with little to no color features.

Table IV. Surface reconstruction accuracy on the synthetic ICL-NUIM dataset by [Handa et al. 2014].

	kt0	kt1	kt2	kt3
DVO SLAM	3.2cm	6.1cm	11.9cm	5.3cm
RGB-D SLAM	4.4cm	3.2cm	3.1cm	16.7cm
MRSMap	6.1cm	14.0cm	9.8cm	24.8cm
Kintinuous	1.1cm	0.8cm	0.9cm	24.8cm
Elastic Fusion	0.7cm	0.7cm	0.8cm	2.8cm
Redwood (rigid)	2.0cm	2.0cm	1.3cm	2.2cm
Ours	0.5cm	0.6cm	0.7cm	0.8cm

Mean distance of each reconstructed model to the ground truth surface. Note that unlike the other methods listed, Redwood does not use color information and runs offline.

Table V. ATE RMSE on the synthetic augmented ICL-NUIM Dataset by [Choi et al. 2015].

	Living room 1	Living room 2	Office 1	Office 2
Kintinuous	27cm	28cm	19cm	26cm
DVO SLAM	102cm	14cm	11cm	11cm
SUN3D SfM	21cm	23cm	24cm	12cm
Redwood	10cm	13cm	6cm	7cm
Ours	0.6cm	0.5cm	15.3cm	1.4cm

Note that unlike the other methods listed, Redwood does not use color information and runs offline.

C. SIFT PERFORMANCE

We provide an additional performance analysis of our GPU-based SIFT detection and matching strategy, see Table VI. Note that for a 1296×968 image (another Structure sensor color resolution), SIFT detection time increases slightly to ≈ 6.4 ms. We detect ~ 150 features per frame, and ~ 250 per keyframe, for all sequences.

Table VI. SIFT performance for a 640×480 image.

#Features	Time Detect (ms)	Time Match (ms)
150	3.8	0.04
250	4.2	0.07
1000	5.8	0.42

Detection time (including descriptor computation) is per frame, and match time is per image pair (parallelized). On all sequences run, we detect about 150 features per frame, and about 250 per keyframe.


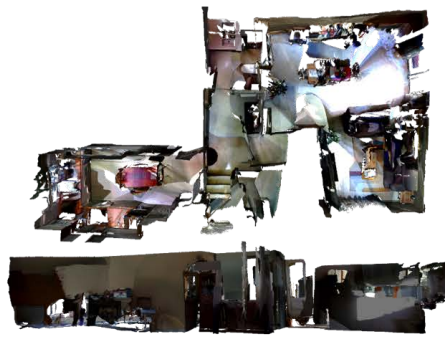
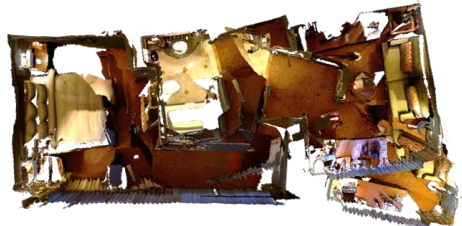
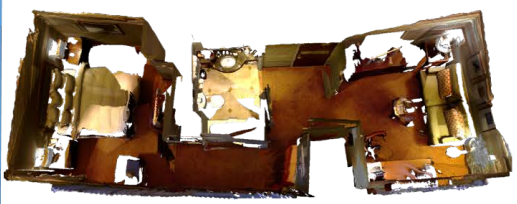






	SUN3Dsfm	Ours
<p>home_at_scan1: 14785 frames 102.6m trajectory</p>		
<p>scan1: 9526 frames 109.4m trajectory</p>		
<p>scan3: 10207 frames 70.1m trajectory</p>		
<p>nips_4: 6269 frames 58.1m trajectory</p>		
<p>hotel_umd_1: 5693 frames 49.6m trajectory</p>		

Fig. 14. Reconstruction results on scenes from the SUN3D dataset [Xiao et al. 2013], using SUN3Dsfm and our approach.

Reconstruction results on SUN3D annotated scenes

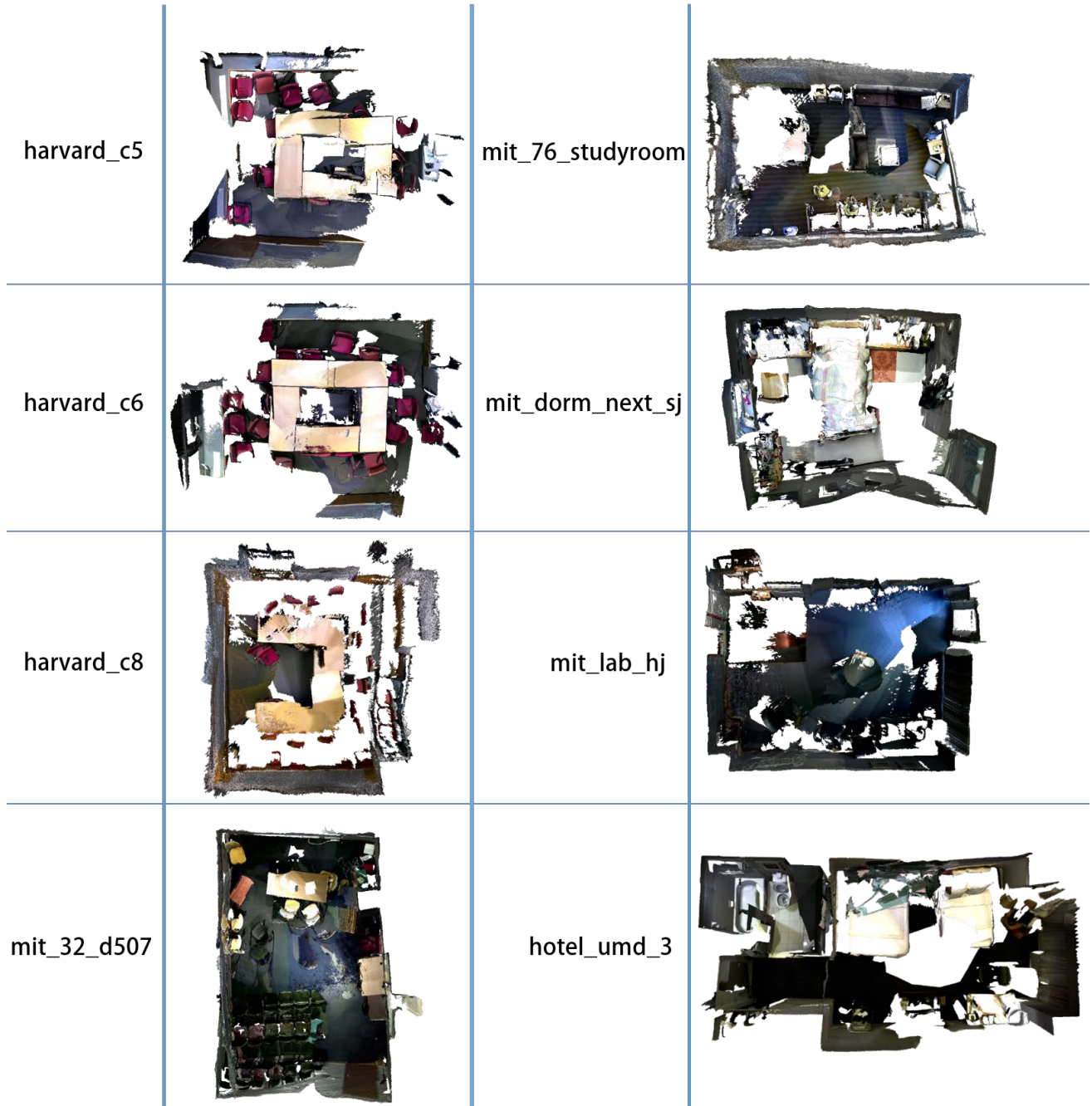


Fig. 15. Reconstruction results on eight scenes from the SUN3D dataset [Xiao et al. 2013], chosen from the *List of Annotated Scenes* (our method is fully automated and does not use any annotations).

Reconstructions of NYU2 Scenes

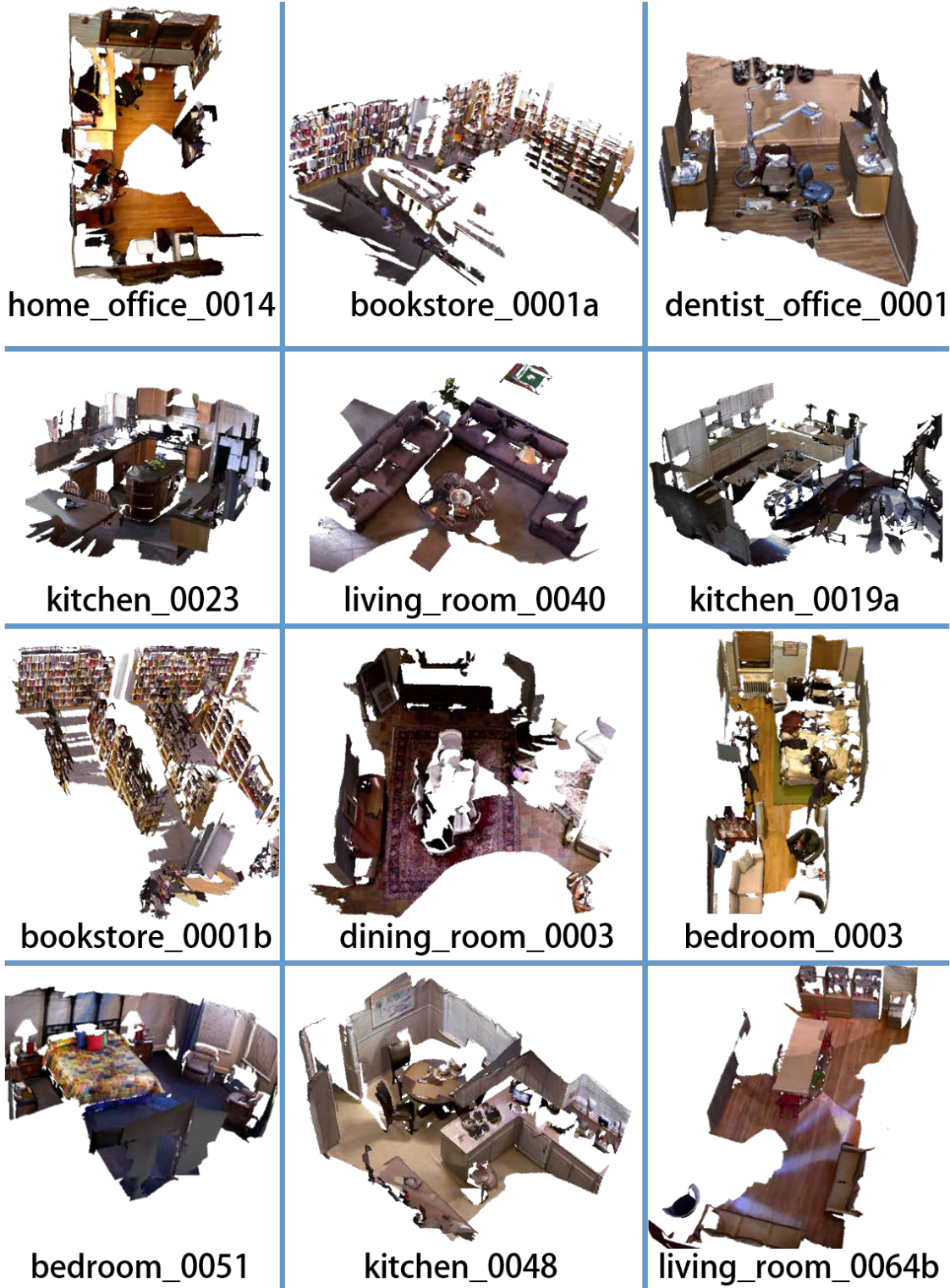


Fig. 16. Reconstructions from the NYU2 dataset [Silberman et al. 2012].