

3D Reconstruction and Texture Optimization Using a Sparse Set of RGB-D Cameras

使用RGBD稀疏set完成3维重建和材质optimization

Wei Li, Xiao Xiao, James K. Hahn

Department of Computer Science, The George Washington University

{gw.liwei, xxfall2012, hahn}@gwu.edu

Abstract

*We contribute a new integrated system designed for high-quality 3D reconstructions. The system consists of a sparse set of commodity **RGB-D** cameras, which allows for fast and accurate scan of objects with multi-view inputs. We propose a robust and efficient tile-based streaming pipeline for geometry reconstruction with **TSDF fusion** which minimizes memory overhead and calculation cost. Our **multi-grid warping method** for texture optimization can address misalignments of both global structures and small details due to the errors in multi-camera registration, optical distortions and imprecise geometries. In addition, we apply **a global color correction method** to reduce color inconsistency among RGB images caused by variations of camera settings. Finally, we demonstrate the effectiveness of our proposed system with detailed experiments of multi-view datasets.*

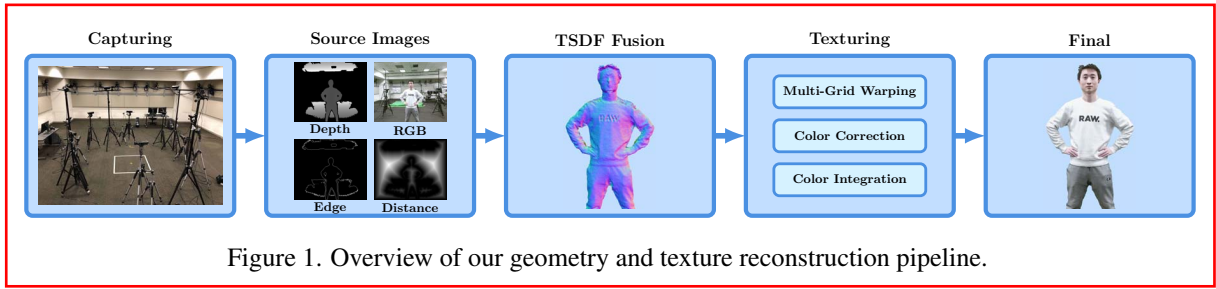
1. Introduction

The ability to capture and reproduce digital versions of real-world objects and scenes is a fundamental and essential task in the fields of multimedia, video games, computer graphics and vision. Recent advances in general purpose computing on GPU (GPGPU), capture technologies (e.g. Microsoft Kinect, 3D scanner) and consumer head-mounted devices (e.g. Microsoft HoloLens, HTC Vive, Oculus Rift) provide the means and the potential for a large variety of applications, such as human body analysis, digitalization of cultural heritage, immersive augmented/virtual reality and many others. As a result, there is an increasing demand for high quality surface reconstruction systems especially with commodity RGB-D sensors.

Some previous researches reconstruct geometries from a sequence of depth images with a single hand-held camera. These methods usually suffer from severe motion blur in the captured RGB images and have synchronization issues between the RGB and the depth cameras. The

camera poses are usually calculated by frame-to-frame or frame-to-model registration that could introduce tracking errors such as drifting. Moreover, the scene needs to be static during the entire scanning process, which is especially difficult for applications like human body scanning [29, 30]. Most of the online reconstruction methods [22, 39, 9, 35, 21, 13, 34, 42, 7] apply truncated signed distance fields (TSDF) to fuse multiple depth images into one volume. Although the sparse data structures of TSDF have been widely researched to reduce memory usage, existing algorithms are still not robust and efficient enough to cull inactive regions for fusion. On the other hand, a visually convincing reconstruction should incorporate both high-quality geometries and textures. Some frameworks simply reconstruct the texture by blending the colors acquired directly from the RGB images [13, 12, 34, 22, 35, 42, 9, 7]. However, since the geometry may not be perfectly reconstructed, and the RGB cameras may not be precisely calibrated and registered, these methods usually produce blurry textures with obvious misalignments. Moreover, most commodity RGB-D devices have no controls over the camera settings (e.g. exposure and white balance), which causes more difficulties in producing color consistent RGB images from different views.

In this paper, we describe how we address the above challenges to achieve high quality reconstruction of photo-realistic textured geometries. Our capturing system consists of a sparse set of commodity RGB-D cameras. The cameras are arranged in a circle surrounding the capture area, which allows for faster and more accurate scan than a single hand-held camera (Fig. 1). We firstly refine the camera poses by global registration to fix general misalignments in the depth and the RGB images. Our tile-based streaming pipeline supports efficient and high resolution GPU-based TSDF fusion. The detection algorithm can robustly cull inactive tiles to reduce unnecessary data transfer and calculation. For reconstructed geometries, we use point-based representation, which allows to extract vertices with adaptive resolution to match the fine details of the RGB images. We then employ a multi-grid warping method to solve the mis-



alignments at multiple scales and address local distortions and edge artifacts at surface boundaries as well. Moreover, we propose a global multi-view color correction method to reduce color variations across all the camera views. Lastly, our weighting scheme in the color integration can ameliorate some of the view-dependent artifacts to further improve the texture quality.

2. Related Work

2.1. TSDF Fusion

Curless *et al.* [12] introduces a volumetric method to incrementally fuse multiple range maps into a single TSDF. The range maps are integrated with weights relative to the uncertainty of measurements, making the method robust against noise. The surface is implicitly represented by the zero crossing level set of TSDF thus the method has no restriction on the topology type of scanned objects. Due to the development of GPGPU and the invention of commodity depth sensors such as Microsoft Kinect, online TSDF based reconstruction has become popular. KinectFusion [22] is one of the real-time surface tracking and reconstruction methods. This method uses dense uniform grids to store TSDF, projective depth difference to estimate signed distance, and point-to-plane ICP to register new measurement with previously fused surfaces. The whole system is implemented on the GPGPU architecture making it possible to perform reconstructions in real time.

However, the underlying data structure of KinectFusion is memory intense, which limits its usage in large-scale and high-detail surface reconstruction. Thus, some researches are focusing on exploiting the spatial sparsity of scanned surfaces to address this problem. Voxel Hashing [35] proposes a spatial hashing scheme to compress and stream voxel chunks. Octomap [21] uses sparse voxel octree to represent the scene in a multi-resolution hierarchy for robot planning and localization. This method also applies a probability model to estimate spatial occupancy for each voxel. Chen *et al.* [9] applies a multi-level hierarchical GPU based data structure to perform dynamic fusion and lossless streaming between GPU and host. In our system we are using tilemaps to subdivide the space into sparse tiles. Although the data structure is not as scalable as octree or hash table, it is sufficient for scanning most indoor

scenes. We will describe our distance and weight metrics in the TSDF fusion and introduce a more robust and efficient culling algorithm to enhance the performance.

2.2. Texture reconstruction

While much attention has been paid to the problem of reconstructing accurate geometries, the reproduction of high-quality textures is equally important in applications that require photo-realistic appearances of the scanned objects. A commonly used approach is to simply project each vertex of the geometry onto the RGB images and compute the vertex color by blending the pixel colors at corresponding image locations. Various sophisticated blending schemes have been proposed by [36, 38, 33, 8]. But since the misalignment of images is not addressed, these methods often result in blurry textures. Instead, some other approaches [25, 40, 44, 10, 27] select one view per face by solving a discrete labeling equation or a Poisson equation to hide visible seams between overlapped images. However, they can still produce tearing artifacts with large misalignments of input images. On the other hand, a number of works propose to address either the inaccuracies in the camera calibration [37, 26, 4, 41, 11] or handle inaccuracies in geometry and optical distortions of the input images [15, 14, 1]. But these methods address each of those core problems separately, and thus, are suboptimal.

More recently, a small number of approaches have addressed the above problems simultaneously. Zhou and Koltun [48] present compelling results on texture optimization by refining rigid camera poses along with non-rigid warpings of RGB frames. However, this method uses a single scale control lattice which may fail to converge to the global minimum. It also lacks regulations on the warping field that may lead to unreasonable deformations. The recent work of Bi *et al.* [5] builds upon the work of Zhou and Koltun [48] by proposing a global patch-based system that synthesizes aligned images from input images. They simultaneously maximize the local similarity between the aligned and input images and ensure the photometric consistency of all the aligned images. This method produces high-quality results in most cases, but is expensive to process, and also not able to preserve the semantic information of the texture when large misalignments are present. In our approach, we address the above shortcomings using a multi-grid warping

method with multiple constraints to regulate the warping field.

However, most of the previous methods [48, 5, 20] assume that the scene is captured with fixed camera settings. Even though their methods present impressive results on texture alignment, the result would suffer from visible color variations when using cameras with automatic parameter settings. To attenuate such artifacts, some methods have been proposed to either compute a leveling texture that alleviates color differences at image borders [43, 3] or search for a linear function that minimizes the color differences between pairs of images [8, 31, 2]. However, the former methods only locally alleviate color differences at overlapping areas of images, while the latter cannot build a very robust color correction function since pixel by pixel alignments between images are not guaranteed. Therefore, we apply a histogram matching algorithm similar to [17, 32] by aligning the cumulative histograms of color correspondences in image pairs, which produces visually better results.

3. Surface Reconstruction

In this section we present our surface reconstruction framework. The depth images acquired from multiple sensors are preprocessed to reduce noises and artifacts, and then globally registered to refine the camera poses. The resulting depth images are fused into TSDF with our distance and weight metrics in a tilemap based streaming pipeline. Finally, point clouds are extracted from TSDF for texture optimization and rendering.

3.1. Depth Image Preprocessing

We preprocess the captured depth images to reduce noises and outliers. We remove the depth samples whose values are out of the reliable measurement range (0.5m to 5m) of the depth cameras. We also discard samples distant from most of the neighboring samples. Moreover, we eliminate small area patches by region growing and further reduce noisy points around surface edges by erosion operation. Finally, depth images are undistorted based on corresponding intrinsic parameters of the cameras and back projected to generate point clouds in each view space.

3.2. Global Registration of Camera Poses

We apply global rigid registration on the generated point clouds P to correct the surface misalignment due to the errors in the camera calibration. Consider two cameras i and j among the K depth cameras. We project the point $\mathbf{p}_i \in P_i$ captured from camera i onto the image plane of camera j . Instead of using projective correspondence, we search in the 5×5 neighboring region around the projection point in P_j and choose the best matching point $\mathbf{p}_{i,j}$ as the correspondence of \mathbf{p}_i . The positions, normals and colors

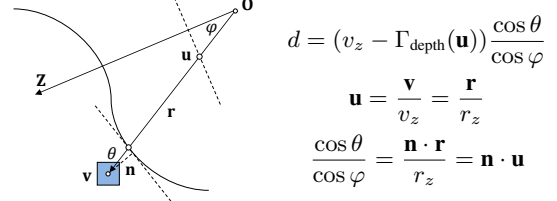


Figure 2. Point-to-plane distance estimation. \mathbf{O} is the camera origin. \mathbf{v} is the voxel center in camera space. \mathbf{u} is the projection of \mathbf{v} on the focal plane. Γ_{depth} fetches the depth value at \mathbf{u} . \mathbf{r} is the viewing direction. ϕ is the angle between \mathbf{r} and the optical axis \mathbf{Z} , while θ is the angle between \mathbf{r} and the local surface normal \mathbf{n} .

of the points are used as features to evaluate the scores of candidate points in the neighboring region. This method results in a better correspondence finding when the surface lacks geometric features or there are misalignments with other surfaces on the edges. For each pair of correspondent points, the residual is measured as their point-to-plane distance [28]:

$$E(\mathbf{T}) = \sum_{j \in K} \sum_{i \in K} \sum_{\mathbf{p}_i \in P_i} \|(T_j M_j M_i^{-1} T_i^{-1} \mathbf{p}_i - \mathbf{p}_{i,j}) \cdot \mathbf{n}_{i,j}\|^2 \quad (1)$$

M are the camera poses from the calibrated extrinsic parameters of the cameras. T are the correction transforms. $M, T \in SE(3)$. Then the registration problem is to solve for the optimal T that minimize the energy function. This results in a dense linear system of $6K$ parameters which is solved by singular value decomposition (SVD) method.

3.3. Fusion with TSDF

After registration, all depth maps are fused into one uniform grid space referred to as the TSDF. Each grid sample of the TSDF at position \mathbf{v} consists of two components: distance $D(\mathbf{v})$ and weight $W(\mathbf{v})$. The distance value measures the distance from the center of the voxel to the closest point on the surface. The reconstructed surface is thus implicitly represented by the zero-crossing level set of the TSDF. It is usually expensive to compute an accurate transform from polygonal meshes to distance fields [23]. We use an alternative fast approach similar to [19] that measures the Euclidean distance from the voxel center to the local tangential plane around the projection point (See Fig. 2). This point-to-plane method only adds a correction term $\cos \theta / \cos \varphi$ to the depth difference but yields a better distance estimation for planar surfaces as it assumes the scanned surfaces to be piecewise linear.

The other component weight is used to incrementally accumulate new measurement $(d_k(\mathbf{v}), w_k(\mathbf{v}))$ from input

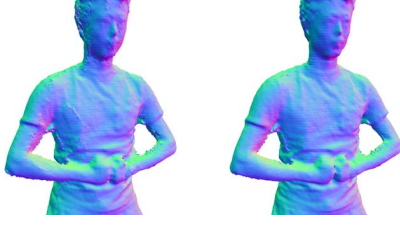


Figure 3. TSDF fusion without (left) and with (right) distance weight w_{edge} . Notice that the visible seams between adjacent scans in the left figure are reduced in the right figure.

depth images with previous measurements in the TSDF:

$$\begin{cases} D_k(\mathbf{v}) = \frac{D_{k-1}(\mathbf{v})W_{k-1}(\mathbf{v}) + d_k(\mathbf{v})w_k(\mathbf{v})}{W_{k-1}(\mathbf{v}) + w_k(\mathbf{v})} \\ W_k(\mathbf{v}) = W_{k-1}(\mathbf{v}) + w_k(\mathbf{v}) \end{cases} \quad (2)$$

The weight is related to the error of the sensor measurement, which has been investigated in several research works [46, 16, 45, 24]. The variance of measurements increases quadratically with respect to the distance. The depth value is also distorted in a ring-like pattern near image borders due to the non-uniform illuminance of the active infrared sensor. Depth inhomogeneity of the time-of-fly (TOF) camera and inappropriate preprocessing of the depth images such as bilateral filtering can introduce non-existing points around the edges of observed surfaces as well. Therefore we design our weighting term based on these metrics as:

$$w = w_{\text{geo}} \cdot w_{\text{border}} \cdot w_{\text{edge}} \quad (3)$$

w_{geo} is the geometric term, which is proportional to the solid angle subtended by the surface patch at the sensor origin. Thus the surface that is closer and facing towards the sensor will contribute more to the final integration. w_{border} is inversely proportional to the square distance of the projection of the point to the optical center of the image to deal with the distortion at image borders and allow smooth transition between overlapped scans. w_{edge} is proportional to the distance from the projection image pixel to the closest discontinuity pixel on the depth image. As depicted in Fig.3, this weight can help reduce artifacts and visible seams near scan edges. The influences of other factors, such as sensor interference, surface materials and environment temperature, are not considered in this paper.

3.4. Streaming Pipeline

A dense volume representation of TSDF is both memory intense and computationally expensive. Therefore, we represent the TSDF in a two-level tilemap grid [18] to exploit the spatial sparsity of scanned objects (Fig. 4). The volume space is subdivided into uniform tiles of size \mathbf{m} . All the allocated tiles of voxels are stored in a large pool. The tile

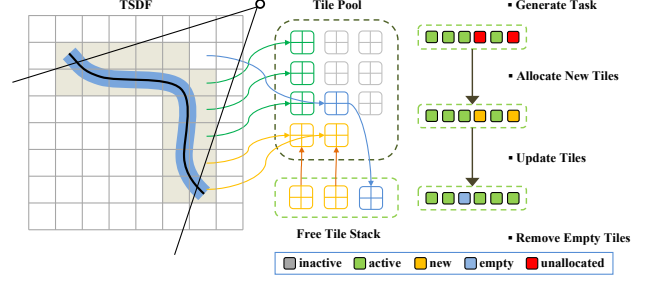


Figure 4. Streaming pipeline of the tilemap based TSDF fusion.

offset in the pool is stored in a separate offset table. The coordinate of a voxel is thus calculated as:

$$\mathbf{p} = \text{OFFSET}(\lfloor \mathbf{v}/\mathbf{m} \rfloor) * \mathbf{m} + \mathbf{v} \text{ MOD } \mathbf{m} \quad (4)$$

where \mathbf{v} is the virtual location of the voxel in the TSDF volume and \mathbf{p} is the corresponding physical location in the tile pool. All the tiles are initialized as free and pushed into a stack.

If a tile is outside the camera frustum, occluded by the scanned surface, or missing valid measurements, the tile is not involved in the fusion phase. Thus we launch a detection kernel to cull these inactive tiles to avoid unnecessary data transfer and calculation. We construct an axis aligned bounding box (AABB) for each tile to perform intersection tests with the view frustum and the scanned surfaces. The size of the AABB should be adjusted to take the width of the TSDF band into consideration. If a tile is not culled by the view frustum, we project the AABB of the tile onto the image plane. We calculate the minimum and the maximum depth values within the projection region as the near and far clipping planes. A tile is culled if the tile is behind the far clipping plane or the tile is empty and in front of the near clipping plane (See Fig. 5). Notice that the empty space in front of the measured depth implies that there is no other surface in the way. Thus, the non-empty tiles in front of the near clipping plane may require integration to fix errors and thus cannot be culled.

The remaining active tiles are pushed into a task list and examined to determine if they have already been allocated in the tile pool. If not, a new tile is popped from the free tile stack and assigned to the unallocated grid. When the allocation is completed, the fusion kernel is launched to update the TSDF value of the scheduled tiles in the task list with the input depth images. Another kernel examines the distance values and removes the tiles that no longer intersect the surface band back to the free tile stack.

3.5. Point Cloud Extraction

A point cloud is extracted from the fused TSDF to prepare for texture optimization and rendering. Each TSDF

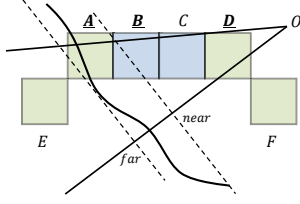


Figure 5. Green tiles are non-empty, while blue tiles are empty. Non-empty tiles F and E are either out of the view frustum or behind the far clipping plane, thus they are both culled. Empty tile C is in front of the near clipping plane, therefore it is culled as well. The other tiles A , B and D are active.

cell is subdivided until the projection of the cell on the color image is smaller than a pixel. Point splats are then extracted at the center of the solid cells on the isosurface. The solid cells are those with opposite signs of distance values at corners. The splat normal is calculated as the gradient of the TSDF by central difference. The splat radius is set to at least $\sqrt{3}/2$ times the cell size to cover the entire cell so that the splat can blend with neighboring splats during rendering. Since the reconstructed surface may be different from the actual scanned object, we generate a new depth image by ray casting the TSDF for each RGB camera view. The new depth image is then used in the texture optimization phase to determine the visibility of a point to the camera view. A GPU-based multi-pass surface splatting technique [6] is implemented to achieve a high-quality rendering result.

4. Texture Optimization

This section explains in detail the texture optimization process aimed at generating a photometrically consistent texture from multi-view RGB images. The process involves three procedures: color integration, multi-grid warping and multi-view color correction. The multi-grid warping fixes the misalignments between the reconstructed geometries and the RGB images, while the color correction improves the color consistency across the RGB images. Color integration is applied in both of the above procedures to produce the proxy of vertex colors for optimizations.

4.1. Color Integration

Given the reconstructed point cloud P with a set of multi-view RGB images $\{I_k\}$, our goal is to compute an optimal color $C(\mathbf{p})$ for each vertex $\mathbf{p} \in P$. We first project each vertex \mathbf{p} onto the generated depth image associated with each RGB image to determine the visibility. If the depth value of the vertex is not greater than the corresponding value in the depth image plus a threshold, the vertex is considered visible to the corresponding RGB image. Then we compute an optimal color $C(\mathbf{p})$ for each vertex as the



Figure 6. Left: optimization only at the finest scale. Middle: our multi-grid approach. The thin line patterns are aligned. Right: an illustration of the control points from a coarser scale (red dots) and the next finer scale (green dots).

weighted average of the pixel colors from all the visible images $\mathcal{V}(\mathbf{p})$:

$$C(\mathbf{p}) = \left(\sum_{k \in \mathcal{V}(\mathbf{p})} w_k(\mathbf{p}) \Gamma_{I_k}(\mathbf{u}_k) \right) / \left(\sum_{k \in \mathcal{V}(\mathbf{p})} w_k(\mathbf{p}) \right) \quad (5)$$

where \mathbf{u}_k is the projection coordinate of \mathbf{p} on the image plane of I_k . $\Gamma_{I_k}(\mathbf{u}_k)$ is the bilinear interpolated pixel color at \mathbf{u}_k . w_k is the weighting term that specifies the contribution of the pixel color at \mathbf{u}_k to the integrated color $C(\mathbf{p})$. A good weighting function can significantly improve the visual quality of the resulting texture reconstruction. Here we reuse the weighting metrics introduced in Sec. 3.3. We discard observations in $\mathcal{V}(\mathbf{p})$ that have less weights or saturated colors to avoid false, clamped or over smoothed values.

4.2. Multi-Grid Warping

Although the global registration of camera poses can solve most of the misalignments among the scans, the result from direct color integration still suffers from blurring artifacts caused by imprecise geometry reconstruction and camera registration. Therefore, we propose a multi-grid warping method to correct these misalignments. Similar to Zhou and Koltun’s method [48], we define a set of control points $\mathbf{Q} : \{\mathbf{q} \in \mathbb{R}^2\}$ and corresponding warping vectors $\mathbf{F} : \{\mathbf{f} \in \mathbb{R}^2\}$ organized in a uniform grid over each image plane (Fig. 6). For any in-between position \mathbf{u} on the image plane, the warping function F is calculated by bilinear interpolating of the warping vectors from its four closest control points.

We first formulate the data term of our energy function as to maximize the agreement of optimal vertex colors $\mathbf{C} : \{C(\mathbf{p})\}$ and corresponding pixel colors in associated images $\mathcal{V}(\mathbf{p})$:

$$E_{\text{data}}(\mathbf{C}, \mathbf{F}) = \sum_{\mathbf{p} \in P} \sum_{k \in \mathcal{V}(\mathbf{p})} \|\Gamma_{I_k}(\mathbf{u}_k + F_k(\mathbf{u}_k)) - C(\mathbf{p})\|^2 \quad (6)$$

However, with sparse input of camera views, not all surfaces of the scanned object can be covered by sufficient samples, which makes the misassignments at the surface edges may

not be simply recovered by the data term. Therefore, we introduce a new term E_{edge} in our energy function:

$$E_{\text{edge}}(\mathbf{F}) = \sum_{k \in K} \sum_{\mathbf{u} \in EG_k^{\text{depth}}} \|\Gamma_{DS^{\text{rgb}}_k}(\mathbf{u}_k + F_k(\mathbf{u}_k))\|^2 \quad (7)$$

EG^{depth} is the set of discontinuity points on the edges of the depth image. DS^{rgb} is the distance map converted from the edge map of the RGB image. This term attempts to align the edges detected in the depth images to those detected in the associated RGB images. This can significantly improve the local texture alignments at surface boundaries (not just between foreground and background objects) and affect the warping vectors within the surface by the regulation terms. This term may introduce artifacts when the edges of RGB images are produced by other sources of discontinuities such as textures or shadows. But since we limit this term within a narrow band around the depth edges, in practice the artifacts are kept to a minimum and hardly noticeable after color integration.

As the warping model above could easily produce some unreasonable deformations, we add another two energy terms E_{smooth} and E_{reg} to regulate the warping vectors. E_{smooth} is applied to ensure that the deformation is consistent across the warping field. Specifically, the warping correction applied to a control point \mathbf{q} should match those applied from its neighboring control points $\mathcal{N}(\mathbf{q})$:

$$E_{\text{smooth}}(\mathbf{F}) = \sum_{k \in K} \sum_{i \in Q_k} \sum_{j \in \mathcal{N}(\mathbf{q}_{k,i})} w_{i,j} \|\mathbf{f}_{k,i} - \mathbf{f}_{k,j}\|^2 \quad (8)$$

where $w_{i,j}$ is a weighting term inversely proportional to the distance between control points. E_{reg} is an L^2 regularizer on the magnitude of the warping vectors $\mathbf{f}_{k,i}$ to keep the deformed control points close to their original positions:

$$E_{\text{reg}}(\mathbf{F}) = \sum_{k \in K} \sum_{i \in Q_k} \|\mathbf{f}_{k,i}\|^2 \quad (9)$$

Finally, our complete energy function is:

$$E(\mathbf{C}, \mathbf{F}) = \lambda_{\text{data}} E_{\text{data}}(\mathbf{C}, \mathbf{F}) + \lambda_{\text{edge}} E_{\text{edge}}(\mathbf{F}) + \lambda_{\text{smooth}} E_{\text{smooth}}(\mathbf{F}) + \lambda_{\text{reg}} E_{\text{reg}}(\mathbf{F}) \quad (10)$$

where λ_{data} , λ_{edge} , λ_{smooth} , λ_{reg} are the weights to balance the strength of each term. The objective $E(\mathbf{C}, \mathbf{F})$ is a non-linear least-squares function. We solve \mathbf{C} and \mathbf{F} alternatively. That is, we first fix the warping function \mathbf{F} and solve \mathbf{C} for each vertex by the color integration, and then fix \mathbf{C} and solve the warping function \mathbf{F} by Levenberg–Marquardt method for each RGB image. The optimization problem results in a large sparse linear system which is efficiently solved by preconditioned conjugate gradient (PCG) method. This process is repeated until the solution converges.

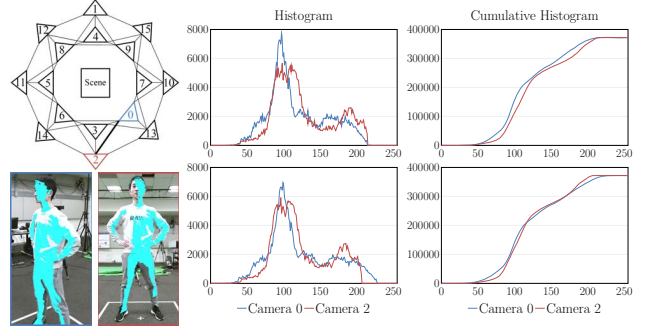


Figure 7. Global multi-view color correction. The edges in the top left graph depict the image pairs with common corresponding pixels. The two figures at the bottom left are examples of the masks of the common pixels in images 0 and 2. The graphs on the right illustrate the intensity distribution of camera 0 and 2 before (first row) and after color correction (second row).

The major issue of this method is how to choose an appropriate size of the control grids. Small grid size can lead to a large sparse linear system that is expensive to solve and tends to be trapped into local minima (See Fig. 6). Large grid size may have difficulties in recovering fine details of the texture. Therefore, we propose a multi-grid method to optimize the warping field \mathbf{F} in a coarse-to-fine manner. Specifically, we start with a coarser resolution of control points and solve for the \mathbf{C} and \mathbf{F} as mentioned before. Then we upsample the warping field from lower resolution by bilinear interpolation (See Fig. 6) and solve for the \mathbf{C} and \mathbf{F} in a finer scale. We repeat these processes until reaching the finest level of control grids. This multi-grid method is more robust and converges in fewer iterations than single scale approaches.

4.3. Color Correction

The RGB cameras of some commodity sensors, such as Microsoft Kinect V2, only allow automatic color and exposure controls, which makes it difficult to produce color consistent images across different camera views especially under non-uniform lighting and background conditions. The resulting texture often reveals visible seams in overlapped regions of different scans (Fig. 8b). The weighting scheme in Eq. 5 can only hide the artifacts by blending the overlapped color samples, but cannot fundamentally solve the problem. Therefore, we present a global multi-view color correction method that optimizes per-camera parameters to set all RGB images in a common color reference.

We begin by building color correspondences between each pair of RGB images. We project each vertex of the geometry onto the depth images associated with the RGB images to perform visibility tests as previously described in

Sec. 4.1. If one vertex can be seen in both camera views, we project the vertex again onto the RGB images to obtain a pair of color values and add them into the list of color correspondences. In order to build a robust color correction function, we discard the color values near the extremes which are less reliable. In addition, if two RGB images do not have enough correspondence points, they are not involved in the optimization phase. All the remaining camera pairs thus form a graph structure G (Fig. 7).

Since small errors in the geometry may result in a large difference in color mapping, pixel by pixel alignment is not guaranteed. Our color correction algorithm is thus based on the distribution of colors. Similar to the method in [32], we build histograms for the correspondent colors for each RGB camera and calculate corresponding cumulated histograms (Fig. 7). Color consistency is solved by optimally aligning the quantile positions of the cumulated histograms to a common one from the reference camera. We define two parameters, gain g and offset o , for each camera to fix color discrepancies among RGB images. We solve the problem as a graph-structured optimization over all the image pairs, which involves a global minimization of a Linear Programming problem:

$$\arg \min_{g \geq 0, g, o \in \mathbb{R}} \max_{(i,j) \in G} |(g_i R_{ij}^l + o_i) - (g_j R_{ji}^l + o_j)| \quad (11)$$

$g_{\text{ref}} = 1 \wedge o_{\text{ref}} = 0$

R_{ij}^l corresponds to the l^{th} quantile of the cumulated histogram of common pixels of image i with image j . ref is the index of the chosen reference image. Here we only use a set of quantile positions instead of the whole set of colors to ensure scalability. We perform the optimization for each color channel (R, G, B) independently.

5. Results

System Our capturing system consists of one master computer, 3 node computers and 16 Kinect v2 cameras. The node computers are connected to the master computer via high speed Ethernet cables. Each computer connects up to 4 Kinects. The Kinects are arranged circularly covering a scanning area of $5\text{m} \times 5\text{m}$ (Fig. 1). Each node computer records depth and RGB images from connected cameras to local disks during capturing and streams the recorded image sequences to the master computer for geometry and texture reconstruction.

TSDF fusion Table 1 demonstrates how our tile-based streaming pipeline can help to reduce TSDF storage and calculation. We fuse the depth images of our benchmark models into a volume at 512^3 and 1024^3 resolution with a tile size of 8 using our culling algorithm described in Sec. 3.4. The second column is the percentage of tiles allocated in the tile pool among all the tiles in the volume. The third

column is the percentage of active tiles processed per depth image on the average. In contrast, KinectFuion [22] needs to process all voxels for each depth image. Voxel hashing [35] performs ray casting through the truncation region around the depth samples to detect active tiles. This method may miss potential tiles on distant surfaces since it only uses rays rather than frustums through the depth samples. Scalable volumetric surface reconstruction [9] applies a similar approach to ours that projects the bounding volumes of voxels to the depth images to detect active voxels. But for the tiles in front of the truncation region, the method may mark them as free space which has the possibility to ruin valid reconstructions by noisy measurements.

Model	Allocated		Active (avg.)	
	512	1024	512	1024
sub-49-08	23.7%	16.5%	9.4 %	5.9%
sub-42-39	32.3%	22.6%	12.1 %	7.7%
sub-56-28	22.2%	14.5%	8.0%	4.7%

Table 1. TSDF storage

Optimization performance We have tested both CPU and GPU implementations for the non-rigid texture warping optimization. The calculation of $J^T J$ and $J^T r$ of the Gaussian-Newton method costs 44ms per iteration for each camera on average at the grid size of 8. In contrast, our block based GPU implementation costs 11ms. The linear system solver is implemented by the Preconditioned Conjugate Gradient (PCG) method which is fixed to 10 iterations. The sparse matrix based CPU version takes 28ms while the GPU version takes 7ms.

Texture quality Fig. 8 presents qualitative results of our texture optimization. Fig. 8a is using naïvely blending for comparison. As can be seen, the results suffer from blurring and ghosting artifacts caused by the misalignments between RGB images and geometries. In Fig. 8b, with our multi-grid warping method, most of the misalignments are corrected and the resulting texture quality is significantly improved. However, visible seams are still present along the transitions of adjacent camera views due to different camera settings. In Fig. 8c, we apply the color correction method which greatly reduces the color variation. In order to clearly show the difference, we render Fig. 8b and Fig. 8c by only assigning each vertex the color from the best view. In Fig. 8d, we apply the weighted color integration. The artifacts between scans completely disappear from the result.

We also evaluate the effects of enforcing our warping energy function with the edge term and the smoothness term (Fig. 9). By optimizing with the original warping energy function described by Zhou and Koltun [48], which only

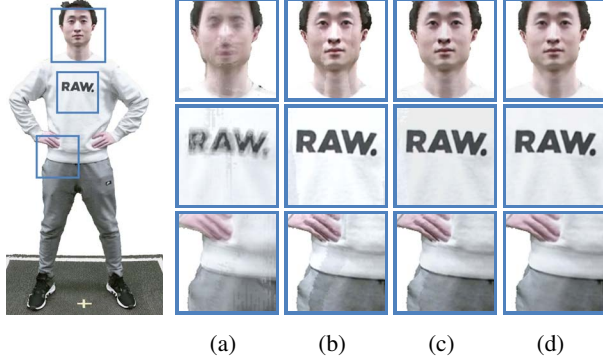


Figure 8. Results of texture optimization. (a) naïve averaging; (b) multi-grid warping; (c) color correction; (d) final result with weighted color integration. (b) and (c) are only using the best view to show the visible seams between camera views.

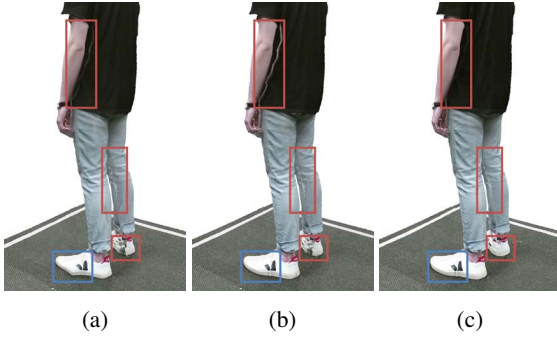


Figure 9. The effects of the smoothness term and the edge term in our warping energy function. From left to right: without smoothness and edge terms; with smoothness term; and with both smoothness and edge terms.

has the data term and the regularization term, the resulting texture can produce visible edge artifacts at surface boundaries (red boxes in Fig. 9a) and unreasonable deformations (blue box in Fig. 9a). By adding the smoothness term to the optimization, the deformations are constrained by neighboring warping field, but the edge artifacts still exist (Fig. 9b). By adding the smoothness term along with the edge term, the alignments on the edges are solved as well.

Fig. 10 compares our texture optimization results against those of the other methods. While volumetric blending [34, 22, 35, 7] produces blurry colors, best labeling method by Waechter *et al.* [44] yields sharper textures. However, the texture consistency may be violated when adjacent vertices are assigned colors from RGB images with large misalignment (notice the tearing artifacts on the face of the subject). Zhou and Koltun’s method [48] can correct most of the misalignments, but with fixed size of control grids the optimization may fail to converge to the global minimum,



Figure 10. A comparison with other optimization algorithms.

which results in undesirable misalignment. Moreover, as described previously, since the warping function has no regulations from the surface boundaries or the neighboring warping field, the resulting texture mapping may have unreasonable distortions and edge artifacts. The problem becomes even worse when there is only a sparse set of image inputs. In contrast, our multi-grid method addresses all these issues properly and achieves a higher quality of texture reconstruction.

6. Conclusion

In this paper, we have presented an integrated system to reconstruct high-quality textured geometries by using a sparse set of RGB-D cameras. Our system is based on low-cost commodity devices, which is flexible and accessible for the creation of VR and AR content for various applications. Our tile-based stream pipeline can perform efficient TSDF fusion with less calculation and memory overhead. Our high-quality texture reconstruction depends critically on the global camera pose registration and multi-grid image warping to solve texture misalignments. In addition, our global multi-view color correction method can further ameliorate visible seams between adjacent camera views due to variations of camera settings and significantly improve the color consistency of the texture.

Acknowledgements This work was supported in part by NSF Grant CNS-1337722 and NIH Grant 1R21HL124443.

References

- [1] E. Aganj, P. Monasse, and R. Keriven. Multi-view texturing of imprecise mesh. In *Asian Conference on Computer Vision*, pages 468–476. Springer, 2009.
- [2] D. S. Alexiadis, A. Chatzitofis, N. Zioulis, O. Zoidi, G. Louizis, D. Zarpalas, and P. Daras. An integrated platform for live 3d human reconstruction and motion capturing. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(4):798–813, 2017.
- [3] M. Arian, R. Preiner, C. Scheiblauer, S. Jeschke, and M. Wimmer. Large-scale point-cloud visualization through localized textured surface reconstruction. *IEEE transactions on visualization and computer graphics*, 20(9):1280–1292, 2014.
- [4] F. Bernardini, I. M. Martin, and H. Rushmeier. High-quality texture reconstruction from multiple scans. *IEEE Transactions on Visualization and Computer Graphics*, 7(4):318–332, 2001.
- [5] S. Bi, N. K. Kalantari, and R. Ramamoorthi. Patch-based optimization for image-based texture mapping. *ACM Transactions on Graphics (Proceedings of SIGGRAPH 2017)*, 36(4), 2017.
- [6] M. Botsch, A. Hornung, M. Zwicker, and L. Kobbelt. High-quality surface splatting on today’s gpus. In *Point-Based Graphics, 2005. Eurographics/IEEE VGTC Symposium Proceedings*, pages 17–141. IEEE, 2005.
- [7] E. Bylow, J. Sturm, C. Kerl, F. Kahl, and D. Cremers. Real-time camera tracking and 3d reconstruction using signed distance functions. In *Robotics: Science and Systems*, volume 2, 2013.
- [8] M. Callieri, P. Cignoni, M. Corsini, and R. Scopigno. Masked photo blending: Mapping dense photographic data set on high-resolution sampled 3d models. *Computers Graphics*, 32(4):464–473, 2008.
- [9] J. Chen, D. Bautembach, and S. Izadi. Scalable real-time volumetric surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(4):113, 2013.
- [10] M. Chuang, L. Luo, B. J. Brown, S. Rusinkiewicz, and M. Kazhdan. Estimating the laplace-beltrami operator by restricting 3d functions. *Computer Graphics Forum*, 28(5):1475–1484, Jul 1, 2009.
- [11] M. Corsini, M. Dellepiane, F. Ganovelli, R. Gherardi, A. Fusiello, and R. Scopigno. Fully automatic registration of image sets on approximate geometry. *International journal of computer vision*, 102(1-3):91–111, 2013.
- [12] B. Curless and M. Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312. ACM, 1996.
- [13] A. Dai, M. Nießner, M. Zollhöfer, S. Izadi, and C. Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (TOG)*, 36(3):24, 2017.
- [14] M. Dellepiane, R. Marroquim, M. Callieri, P. Cignoni, and R. Scopigno. Flow-based local optimization for image-to-geometry projection. *IEEE Transactions on Visualization and Computer Graphics*, 18(3):463–474, 2012.
- [15] M. Eisemann, B. De Decker, M. Magnor, P. Bekaert, E. De Aguiar, N. Ahmed, C. Theobalt, and A. Sellent. Floating textures. In *Computer graphics forum*, volume 27, pages 409–418. Wiley Online Library, 2008.
- [16] P. Fankhauser, M. Bloesch, D. Rodriguez, R. Kaestner, M. Hutter, and R. Siegwart. Kinect v2 for mobile robot navigation: Evaluation and modeling. In *Advanced Robotics (ICAR), 2015 International Conference on*, pages 388–394. IEEE, 2015.
- [17] U. Fecker, M. Barkowsky, and A. Kaup. Histogram-based prefiltering for luminance and chrominance compensation of multiview video. *IEEE Transactions on Circuits and Systems for Video Technology*, 18(9):1258–1267, 2008.
- [18] T. Fogal, A. Schiewe, and J. Krüger. An analysis of scalable gpu-based ray-guided volume rendering. In *Large-Scale Data Analysis and Visualization (LDAV), 2013 IEEE Symposium on*, pages 43–51. IEEE, 2013.
- [19] S. F. Frisken and R. N. Perry. Efficient estimation of 3d euclidean distance fields from 2d range images. In *Volume Visualization and Graphics, 2002. Proceedings. IEEE/ACM SIGGRAPH Symposium on*, pages 81–88. IEEE, 2002.
- [20] Y. Fu, Q. Yan, L. Yang, J. Liao, and C. Xiao. Texture mapping for 3d reconstruction with rgb-d sensor. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4645–4653, 2018.
- [21] A. Hornung, K. M. Wurm, M. Bennewitz, C. Stachniss, and W. Burgard. Octomap: An efficient probabilistic 3d mapping framework based on octrees. *Autonomous Robots*, 34(3):189–206, 2013.
- [22] S. Izadi, D. Kim, O. Hilliges, D. Molyneaux, R. Newcombe, P. Kohli, J. Shotton, S. Hodges, D. Freeman, A. Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568. ACM, 2011.
- [23] M. W. Jones. The production of volume data from triangular meshes using voxelisation. In *Computer Graphics Forum*, volume 15, pages 311–318. Wiley Online Library, 1996.
- [24] E. Lachat, H. Macher, M. Mittet, T. Landes, and P. Grussenmeyer. First experiences with kinect v2 sensor for close range 3d modelling. *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 40(5):93, 2015.
- [25] V. Lempitsky and D. Ivanov. Seamless mosaicing of image-based texture maps. pages 1–6. IEEE, 2007.
- [26] H. P. Lensch, W. Heidrich, and H.-P. Seidel. A silhouette-based algorithm for texture registration and stitching. *Graphical Models*, 63(4):245–262, 2001.
- [27] H. Li, E. Vouga, A. Gudym, L. Luo, J. T. Barron, and G. Gusev. 3d self-portraits. *ACM Transactions on Graphics (TOG)*, 32(6):187, 2013.
- [28] K.-L. Low. Linear least-squares optimization for point-to-plane icp surface registration. *Chapel Hill, University of North Carolina*, 4, 2004.
- [29] Y. Lu, S. McQuade, and J. K. Hahn. 3d shape-based body composition prediction model using machine learning. In

- 2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), pages 3999–4002. IEEE, 2018.
- [30] Y. Lu, S. Zhao, N. Younes, and J. K. Hahn. Accurate non-rigid 3d human body surface reconstruction using commodity depth sensors. *Computer Animation and Virtual Worlds*, 29(5):e1807, 2018.
 - [31] A. Maimone and H. Fuchs. Encumbrance-free telepresence system with real-time 3d capture and display using commodity depth cameras. pages 137–146, 2011.
 - [32] P. Moulon, B. Dussan, and P. Monasse. Global multiple-view color consistency. In *CVMP*, pages to–appear, 2013.
 - [33] P. J. Neugebauer and K. Klein. Texturing 3d models of real world objects from multiple unregistered photographic views. *Computer Graphics Forum*, 18(3):245–256, Sep 1999.
 - [34] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*, pages 127–136. IEEE, 2011.
 - [35] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 32(6):169, 2013.
 - [36] E. Ofek, E. Shilat, A. Rappoport, and M. Werman. Multiresolution textures from image sequences. *IEEE Comput. Graph. Appl.*, 17(2):18–29, Mar. 1997.
 - [37] F. Pighin, J. Hecker, D. Lischinski, R. Szeliski, and D. Salesin. Synthesizing realistic facial expressions from photographs. SIGGRAPH '98, pages 75–84. ACM, Jul 24, 1998.
 - [38] K. Pulli, H. Abi-Rached, T. Duchamp, L. G. Shapiro, and W. Stuetzle. Acquisition and visualization of colored 3d objects. In *Proceedings. Fourteenth International Conference on Pattern Recognition (Cat. No.98EX170)*, volume 1, pages 11–15 vol.1, Aug 1998.
 - [39] H. Roth and M. Vona. Moving volume kinectfusion. In *BMVC*, volume 20, pages 1–11, 2012.
 - [40] S. N. Sinha, D. Steedly, R. Szeliski, M. Agrawala, and M. Pollefeys. Interactive 3d architectural modeling from unordered photo collections. *ACM Transactions on Graphics*, 27(5):1, Dec 1, 2008.
 - [41] I. Stamos and P. K. Allen. Geometry and texture recovery of scenes of large scale. *Computer Vision and Image Understanding*, 88(2):94–118, 2002.
 - [42] F. Steinbrucker, C. Kerl, and D. Cremers. Large-scale multi-resolution surface reconstruction from rgb-d sequences. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3264–3271, 2013.
 - [43] L. Velho and J. Sossai Jr. Projective texture atlas construction for 3d photography. *The Visual Computer*, 23(9-11):621–629, 2007.
 - [44] M. Waechter, N. Moehrle, and M. Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *European Conference on Computer Vision*, pages 836–850. Springer, 2014.
 - [45] O. Wasenmüller and D. Stricker. Comparison of kinect v1 and v2 depth images in terms of accuracy and precision. In *Asian Conference on Computer Vision*, pages 34–45. Springer, 2016.
 - [46] L. Yang, L. Zhang, H. Dong, A. Alelaiwi, and A. El Saddik. Evaluating and improving the depth accuracy of kinect for windows v2. *IEEE Sensors Journal*, 15(8):4275–4285, 2015.
 - [47] L. Yao, Z. Shang, and Y. Naji. Accurate nonrigid 3d human body surface reconstruction using commodity depth sensors. *Computer Animation and Virtual Worlds*, 0(0):e1807. e1807 cav.1807.
 - [48] Q.-Y. Zhou and V. Koltun. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Transactions on Graphics (TOG)*, 33(4):155, 2014.