

Volumetric Capture

Marcel Bruckner

marcel.bruckner@tum.de

Kevin Bein

kevin.bein@tum.de

Moiz Sajid

moiz.sajid@tum.de

Technische Universität München

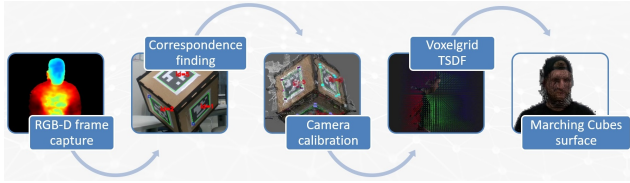


Figure 1: Reconstruction pipeline used in this project: RGB-D Frame Capture, Correspondence Finding, Camera Calibration, Voxelgrid TSDF, Marching Cubes

Abstract

Our goal in this project is to get a real time mesh reconstruction of dynamic scenes using a multi view camera capture setup. We use three Intel RealSense D415 depth cameras for RGB-D data capturing which are calibrated using known correspondences from markers. The data is fused into a voxelgrid representing the implicit surface and following a marching cubes algorithm is performed for the final surface extraction. Both qualitative and quantitative results are presented.

1. Introduction

Volumetric Capture is an extensively researched topic where the goal is to get an accurate 3D reconstruction of dynamic scenes in real time. The first work in the area was done by Cureless and Levoy [1] in which they integrated data from one range image into a cumulative weighted signed distance function. KinectFusion [8] was another milestone in this area which proposed a real time mapping system of indoor scenes using a single depth camera. DynamicFusion [7] handles non-rigid scenes and reconstruction over time. Our work mostly focuses on reconstructing dynamic scenes on a per frame basis.

One application area of our project is Holoportation [13] where 3D reconstructed models are used for VR interactions between two persons over huge distances.

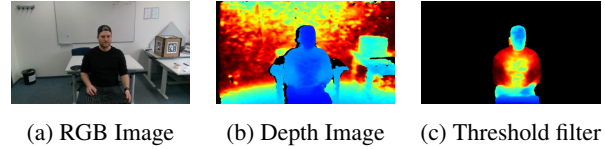


Figure 2: The color and depth images captured by the cameras

2. Reconstruction Pipeline

To perform the real time 3D reconstruction we implemented the pipeline which is exemplary displayed in Figure 1 and described in the following section.

All tests were performed on a Laptop with an Intel® Core™ i7-6700HQ CPU @ 2.60GHz [3], 40GB of DDR4 RAM and a Nvidia GeForce GTX 970M [9].

2.1. RGB-D Frame Capture

We use three Intel® RealSense™ Depth Camera D415 [4] for our capture setup which work on the concept of active stereo depth. Two cameras are placed on the sides of the camera and the depth is calculated from a triangulation using the displacement. An infrared pattern is projected into the scene to enhance the details.

We process the depth images using a threshold filter which discards pixels that are above or below a threshold set during runtime. We also tried different hole-filling, spatial and edge enhancing filters that did not improve the results and are not used in the final reconstruction.

Figure 2 displays the raw RGB-D images and the filtered depth image.

We tried to extend the three camera setup with a fourth one but were not able to integrate it due to bandwidth limitations. Furthermore, we had to limit the color and depth resolutions as well as framerates to enable three cameras and to allow a processing of the data in real time. The resolution we used for the color and depth stream is 848x480 px at 30 frames per second. [6] gives more detail about the limitations of the cameras in a multi camera setup.

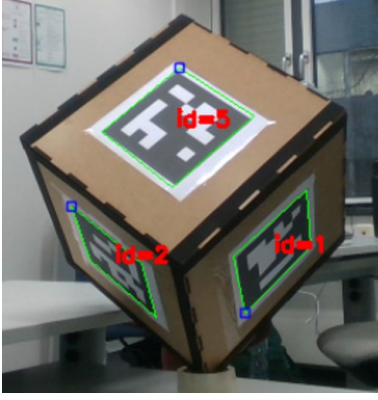


Figure 3: Marker cube with the detected markers highlighted in green and their unique identifiers (red).

The RGB-D images from the three cameras are time synchronized [6] what improves the results of the following camera calibration and the final mesh reconstruction.

Intel[©] provides an excellent SDK [5], good support and updates for its API which made it easy to use the cameras.

2.2. Correspondence Finding

We use ArUco markers [10] to find correspondences between the three camera streams. These markers have a unique identifier and can easily be detected in every color image using the OpenCV [12] library. The markers are mounted on a wooden cube box as shown in Figure 3.

By placing the cameras such that each one is looking on one of the cubes corners we can assure that in every color image three markers are visible. This also ensures that between every pair of cameras we have an overlap of two markers for which the eight corner positions are detected sub pixel perfect. This gives us a total of $8 * 3 = 24$ constraints between every pair of cameras which are used in the next step to estimate the relative transformation between the markers.

It is also important to mention here that we tried ChArUco markers [11] which are a mixture of ArUco markers [10] and a chessboard pattern which allow simultaneous detection and pose estimation. However, the ChArUco markers [11] did not give us good pose results and slowed down the calculation (Table 1), so we stick to ArUco markers [10].

2.3. Camera Calibration

We use the known marker positions in the color image and combine them with the depth information to calculate the 3D position of the marker corners. As we ensure to have an overlap of two markers between every pair of cameras

we now have $8 * 3 = 24$ constraints (8 markers each with a x, y and z value). The 3D positions are then fed into two alignment algorithms to find the relative transformations between the cameras and to align the pointclouds based on the results.

2.3.1 Procrustes

The Procrustes algorithm aligns two pointclouds using known correspondences by estimating the relative translation and rotation between the two sets of points.

Translation We calculate the center of gravity for every pointcloud by summing up the individual point positions and dividing by the number of points. This gives us for every set of points the mean point. The translation is then estimated by the difference between the centers of gravity.

$$T_{ij} = C_i - C_j \quad (1)$$

Rotation The rotation between two pointclouds can be estimated by optimizing over the unknown rotation R_{ij} . This is done by minimizing the mean squared error between the set of points X_i and the X_j .

$$\min_{R_{ij}} \|X_i - R_{ij}X_j\|_2^2 \quad (2)$$

Fortunately there exists a closed form solution for the rotation R_{ij} which is based on the singular value decomposition of the matrix $X_i^T X_j$.

$$X_i^T X_j = U \Sigma V^T \quad (3)$$

$$R_{ij} = UV^T \quad (4)$$

Table 1 displays the duration and mean squared error values after the Procrustes alignment. It shows a still fairly high error, so we had to further align the pointclouds.

2.3.2 Point-to-Point Error

In order to further improve the alignment that we get from Procrustes we use a Point-to-Point correspondence error. We optimize for the relative translation T_{ij} and rotation R_{ij} over every pair of camera frames i and j and in each of these for the known correspondences k . We came up with the following energy term:

$$\min_{T_{ij}, R_{ij}} \sum_i \sum_j \sum_k^{frame \ frame \ corres.} \|X_{ik} - T_{ij} R_{ij} * X_{jk}\|_2^2 \quad (5)$$

A comparison of the results from this optimization is given in Table 1. It shows that the mean squared error is much lower using the transformation from the Point-to-Point optimization. On the other side it displays an increase of the

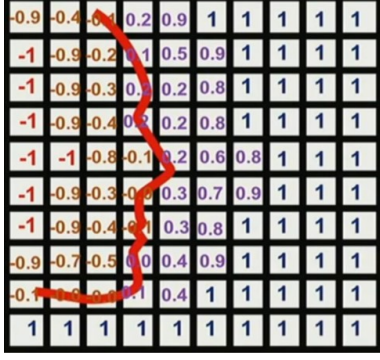


Figure 4: TSDF Representation

duration for the estimation as it is an iterative optimization algorithm, but as the pose estimation is done only once in the beginning the increase is negligible.

2.4. Voxelgrid TSDF

Once we have these aligned pointclouds, we fuse them into the voxelgrid. We project the voxelgrid into the depth frame of the camera and we calculate the difference by subtracting the voxel depth from the actual depth value using Equation (6). For the truncated signed-distance function (tsdf) values we did the weighted averaging using the Equation (7). The tsdf values were truncated between -1 and 1 with points having a negative value inside the surface and positive value outside the surface.

$$tsdf = z_{voxel} - z_{depth} \quad (6)$$

$$tsdf_{i+1} = \frac{tsdf_i * weight}{weight + 1} \quad (7)$$

Given that we were using a large number of gridpoints, computing tsdf values on the CPU was not feasible. So we had to move our tsdf voxelgrid implementation onto the GPU in order to achieve real time results. Since we were already using OpenGL for the rendering pipeline, we decided to use OpenGL Compute Shaders [2] for parallelizing the tsdf voxelgrid implementation. We also had to find a tradeoff between voxel resolution and frame rate. With a lower voxel resolution, the frame rate dropped because more calculations had to be done. Furthermore, finding a good truncation distance was also very crucial for lowering the artefacts. The Table 2 summarizes the time comparisons between different voxel resolutions.

2.5. Surface Extraction

For the iso-surface extraction, we used the Marching Cubes algorithm which converts the implicit surface representation to a polygonal mesh. We then iterated over the

voxelgrid cells to determine the zero crossings and looked up the triangulation in the Marching Cubes table. This gives us the triangulated mesh as a final output. With a lower voxel resolution and more grid points, we were able to get more finer reconstruction. Again we had to utilize the OpenGL Compute Shader in order to achieve real time results because the computations on the CPU were really slow.

Initially, we had voxelgrid points that were not initialized because they were outside the camera view. As a result, we had a huge number of artefacts in the final mesh reconstruction which took a lot of time to figure out. We implemented a two pass OpenGL Compute Shader. In the first pass, we counted the number of triangles that will generated and only allocated space for them. In the second pass, we generate the triangles. We do this because the upper bound for the number of generated triangle is five times the number of voxelgrid points. For example if we have 200x200x200 grid points with a 5 millimeter voxel resolution, we roughly had 40 million number of possible triangles. We only generated a fraction of these triangles after the the triangle counting OpenGL Compute Shader pass. For a voxel resolution of 5mm resolution, the Marching Cubes algorithm took 15 milliseconds which was real time capable. For 3 millimeter resolution, the results were not real time capable however, we were to get fine mesh reconstructions. The Table 2 summarizes the time comparisons between different voxel resolution. The Table 3 summarizes the time comparisons for running Marching Cubes between different voxel resolutions.

3. Results

The results that we get using one camera are really good as can be seen in Figure 5. Figure 6 displays the triangulation results using one camera. However, we get artefacts when we use all three camera streams. The which artefacts which are mostly visible around the edges can be seen in Figure 7. This suggests that the camera pose parameters can be improved further.

Algorithm	Iterations	Duration [ms]	MSE
Procrustes (A)	1	8 - 20	<0.7 - 1.73
Point-to-Point (A)	1 - 20	900 - 1000	<0.1
Procrustes (C)	1	20 - 40	0.25 - 1.29
Point-to-Point (C)	1 - 20	18000 - 22000	<0.1

Table 1: Comparisons between different optimization algorithms in terms of iterations, duration and MSE (Mean Squared Error). A = ArUco, C = ChArUco



Figure 5: Result with using only one camera



Figure 6: Triangulation result with using only one camera

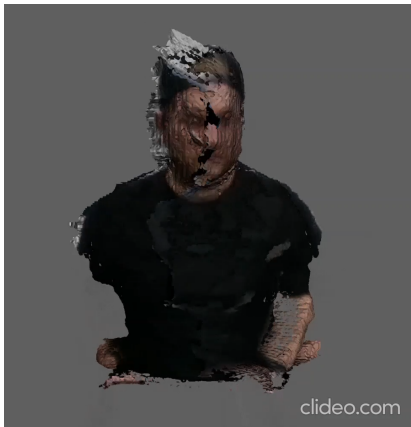


Figure 7: Result with using three cameras

Algorithm	Resolution [mm]	Gridpoints	Duration [ms]
All frames	20	125,000	<1
All frames	10	1,000,000	<1
All frames	5	8,000,000	5
All frames	4	15,625,000	20
All frames	3	37,000,000	50

Table 2: Time comparisons for calculating voxelgrid of different voxel resolutions as well as the number of gridpoints used. The voxelgrid measured 1 meter x 1 meter x 1 meter.

Resolution [mm]	Duration [ms]	Gridpoints
20	<1	7,000 - 10,000
10	<1	40,000 - 45,000
5	15	200,000 - 250,000
4	35	380,000 - 420,000
3	80	600,000 - 800,000

Table 3: Time comparisons for running the Marching Cubes on different voxel resolutions. The number of triangles generated are also displayed.

4. Future Work and Conclusion

Since we optimized for pose parameters for a limited number of markers, one of the further improvements that we can do in camera calibration is to use more correspondences which can be found for example by using Nearest Neighbor Search and then use the Global Alignment techniques like Iterative Closest Points (ICP).

5. Group Member Contributions

References

- [1] B. Curless and M. Levoy. A volumetric method for building complex models from range images. *ACM SIGGRAPH*, 1996.
- [2] A. Gerdelen. It's more fun to compute: An introduction to compute shaders. <http://antongerdelan.net/opengl/compute.html>, 2016. Accessed: 10.02.2020.
- [3] Intel®. Intel® core™ i7-6700HQ processor. <https://ark.intel.com/content/www/de/de/ark/products/88967/intel-core-i7-6700hq-processor-6m-cache-up-to-3-50-g.html>. Accessed: 19.02.2020.
- [4] Intel®. Intel® realsense™ depth camera d415. <https://www.intelrealsense.com/depth-camera-d415/>. Accessed: 19.02.2020.
- [5] Intel®. Intel® realsense™ sdk. <https://www.intelrealsense.com/developers/>. Accessed: 19.02.2020.
- [6] Intel®. Using the intel® realsense™ depth camera d4xx in multi-camera configurations. <https://www.intelrealsense.com/depth-camera-d4xx-in-multi-camera-configurations/>.

[//www.intel.com/content/dam/support/us/en/documents/emerging-technologies/intel-realsense-technology/RealSense_Multiple_Camera_WhitePaper.pdf](https://www.intel.com/content/dam/support/us/en/documents/emerging-technologies/intel-realsense-technology/RealSense_Multiple_Camera_WhitePaper.pdf). Accessed: 19.02.2020.

- [7] R. A. Newcombe, D. Fox, and S. M. Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. *IEEE CVPR*, 2015.
- [8] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. *IEEE ISMAR*, 2011.
- [9] Nvidia. Geforce gtx 970m. <https://www.geforce.com/hardware/notebook-gpus/geforce-gtx-970m>. Accessed: 19.02.2020.
- [10] OpenCV. Detection of aruco markers. https://docs.opencv.org/trunk/d5/dae/tutorial_aruco_detection.html. Accessed: 19.02.2020.
- [11] OpenCV. Detection of charuco corners. https://docs.opencv.org/3.4/df/d4a/tutorial_charuco_detection.html. Accessed: 19.02.2020.
- [12] OpenCV. Opencv. <https://opencv.org/>. Accessed: 19.02.2020.
- [13] S. Orts-Escolano, C. Rhemann, S. Fanello, W.Chang, A. Kowdle, Y. Degtyarev, David Kim, P. Davidson, S. Khamis, M. Dou, V. Tankovich, C. Loop, Q. Cai, P. Chou, S. Men-nicken, J. Valentin, V. Pradeep, S. Wang, S. Bing Kang, P. Kohli, Y. Lutchyn, C. Keskin, and S. Izadi. Holoportation: Virtual 3d teleportation in real-time. *ACM User Interface Software and Technology Symposium (UIST)*, 2016.
- [14] J. Vries. Learnopengl. <https://learnopengl.com/>, 2015. Accessed: 10.02.2020.