



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



3D Depth Sensor on Mobile Devices

Acquisition, Calibration and 3D Reconstruction

Bachelor Thesis

Lukas Bischofberger

lukasbi@ethz.ch

Computer Vision and Geometry Group

Institute for Visual Computing

ETH Zürich

Supervisors:

Dr. Amael Delaunoy

Prof. Dr. Marc Pollefeys

May 15, 2015

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

3D Depth Sensors on Mobile Devices
(Acquisition, Calibration and 3D Reconstruction)

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

Bischofberger

First name(s):

Lukas

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the 'Citation etiquette' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

Signature(s)

Zürich, 24.4.2015



For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.

Acknowledgements

I would like to thank my supervisor Amael Delaunoy for his support and great help I received during the weekly meetings. I also want to thank Professor Marc Pollefeys for offering me the opportunity to write my bachelor thesis at the Computer Vision and Geometry Group.

Abstract

The introduction of Kinect has made depth sensors available to the consumer world. This made it possible for basically everyone to have a 3D scanner at home. Systems like Kinect Fusion would give the possibility to easily get a 3D reconstruction of everything in front of the lens.

The structure sensor by Occipital is the first mobile depth sensor and thus brings the possibility to have a mobile 3D scanner. We want to show that mobile depth data acquisition is possible and easy and also that this data performs well in 3D reconstruction applications.

Therefore a pipeline from data acquisition on a mobile device up to 3D reconstruction on a desktop computer has been developed.

Contents

Acknowledgements	ii
Abstract	iii
1 Introduction	1
1.1 Motivation	1
1.2 Contribution	2
1.3 Related work	2
2 Data Acquisition	4
2.1 Mobile Depth Acquisition	4
2.2 RGB and Depth Camera Calibration	5
2.3 Depth re-projection	11
3 3D Reconstruction	13
3.1 Camera Tracking - SVO	13
3.2 Scale Estimation	14
3.3 Pointcloud	15
3.4 Dense Reconstruction	16
4 App Integration	19
5 Conclusion	22
6 Future work	24
A Code	A-1
A.1 Depth Re-projection	A-1

CHAPTER 1

Introduction

We want to use a state of the art mobile device for depth data acquisition. For this we take advantage of the first mobile depth sensor, the Structure sensor by Occipital¹. Further we will combine the depth sensor with the color camera of the device for texturing a 3D model.

We want to show that this data can then be used to create a 3D reconstruction. For this step we will use camera tracking and motion estimation algorithms and another pipeline for volumetric mapping.

1.1 Motivation

We want to show that with mobile devices like smartphones or a tablets, which today are omnipresent and a fairly low cost sensor we can perform 3D scanning. This scans could be used in numerous applications in daily life.

Scans of rooms can for example be used to measure their dimensions and check for furniture to fit in. This is partly implemented in the IKEA Catalog application, whereas there no real measurements and room structure are available². On the other hand such a scan could also be used for real estate advertisements, e.g. by showing reconstructed interiors depicted with real dimensions. In general every reconstructed object could also be printed by a 3D printer. A last intention is also to show that 3D reconstruction with depth data performs better and faster than such with color images only.

¹<https://occipital.com>

²<https://play.google.com/store/apps/details?id=com.ikea.catalogue.android>

1.2 Contribution

The contributions of our work are the following:

- An iOS application for the depth and color data acquisition which is based on the Structure SDK³⁴. It allows viewing and recording of RGB-D⁵ data as well as camera calibration.
- An offline pipeline which consists of two parts. The first part will convert the raw data recorded on the mobile device to RGB-D pairs using the calibration results. The second part estimates the camera motion and performs a 3D reconstruction.

1.3 Related work

There have been different approaches on reconstruction, ranging between professional solutions, to consumer accessible products up to research projects. They also differ on their technical approach. NextEngine for example uses high precision laser scanning and KinectFusion uses its proper sensor which is a combination of different cameras. The third approach described here solely relies on color images.

1.3.1 Kinect Fusion

Kinect Fusion [4] is an application which provides 3D scanning and model creation based on the Kinect sensor [3]. The Kinect sensor is a horizontal bar which features a RGB camera, a depth sensor and a multi-array microphone [8]. The user can paint a scene with the Kinect camera and simultaneously see, and interact with, a detailed 3D model of the scene. Kinect Fusion runs on a Desktop PC with the sensor attached by wire. The reconstruction and estimation of the camera movement is based on the tracking of the camera through the depth map. Further there exist possibilities to color the constructed models or to apply shading.

1.3.2 Live Metric 3D Reconstruction on Mobile Phones

A paper by Tanskanen et al. [7] proposes a complete on-device approach for 3D reconstruction. It uses a monocular hand-held device without any depth sensor. The reconstruction process is solely based on color images but the camera motion

³Software development kit

⁴<http://structure.io/developers>

⁵RGB color space with an added depth channel

estimation takes advantage of the device own inertial sensors. Furthermore the paper proposes an efficient and accurate multi-resolution scheme for dense stereo matching which takes into account the restricted computational capabilities of the mobile device.

1.3.3 3D Scanner

In the industry products like NextEngine 3 Scanner⁶ have been available for some time. This scanner has a laser for high precision scanning and thus obtains far higher quality models than the approaches described before. It also features export to industry standard applications and allows easy printing of the scans. On the other side special equipment is needed which is considerably more expensive than consumer electronic devices.

1.3.4 Structure Sensor

Our approach is different insofar that we use generally available hardware as the iPad's color camera and an affordable depth sensor. However our data is quite similar to that of [4] with the difference that the data capturing is entirely done on a mobile device. In theory the reconstruction could also be done on the mobile device as described in chapter 5.

⁶<http://www.nextengine.com>

CHAPTER 2

Data Acquisition

In this chapter we give a detailed description on how the data acquisition on a mobile device is performed. In a first step we describe the use of the Structure sensor, its mechanism and how we obtain the data in a meaningful way. Then we describe the relation between the depth data and color images taken from the device camera. In a third step we will show how to perform camera calibration and combine the data to RGB-D image pairs.

2.1 Mobile Depth Acquisition

The Structure sensor is a so called structured-light 3D scanner. It possesses an infrared structured light projector and a corresponding camera system [**structure·dev**]. The projector projects an infrared mesh onto the scene, which the camera records from a slightly different angle. The three-dimensionally shaped surface produces a distorted mesh relative to the viewing camera, which is then able to perform a geometric reconstruction of the surface shape.

The sensor is connected to the device by cable, in our case to the iPad through the lightning port. We use the Structure SDK [**structure·dev**] as an interface between our application and the sensor hardware. This allows us to receive depth frames or infrared images at different scales. Unfortunately the maximum resolution of the depth sensor is only 640x480 pixels. Additionally, the SDK provides the possibility to get synchronized infrared/color or depth/color pairs. The color images are taken by the mobile device's own camera, e.g. the iPad's camera through the AVFoundation framework¹.

Originally we wanted to perform the mobile data acquisition on an Android device. This had the advantage that we could rely on standard libraries like

¹<https://developer.apple.com/library/mac/documentation/AVFoundation/Reference/AVFoundationFramework/index.html>

²http://structure.io/static/occipital/images/v1/getstarted/ipad_camera_centered.jpg



Figure 2.1: Structure sensor on the iPad²

OpenNI³ and OpenCV⁴. The drawback was that applications on Android are written in Java and the C++ libraries are only accessible through the NDK⁵ ⁶ interface. Unfortunately we weren't able to make the structure sensor work on this platform, thus we decided to work on the iOS platform where a SDK for the sensor was provided.

On the data acquisition screen (Figure 2.2) we see the infrared image on the top left, the color image on the top right, the depth map on the bottom left, whereas red/yellow pixels indicate a close location and blue/black pixels represent the far away areas. Finally we can also see the surface normals on the bottom right.

2.2 RGB and Depth Camera Calibration

For the camera calibration we make use of the pinhole camera model which will briefly be explained here. The next section describes how the depth sensor and the device's RGB camera are calibrated. This will provide the camera matrices and lens distortion coefficients for both cameras. After the single camera calibration for both cameras, the third section explains how these results can be used to perform a stereo camera calibration, which finds the extrinsic parameters.

³<http://structure.io/openni>

⁴<http://opencv.org>

⁵Native Development Kit

⁶<https://developer.android.com/tools/sdk/ndk/index.html>

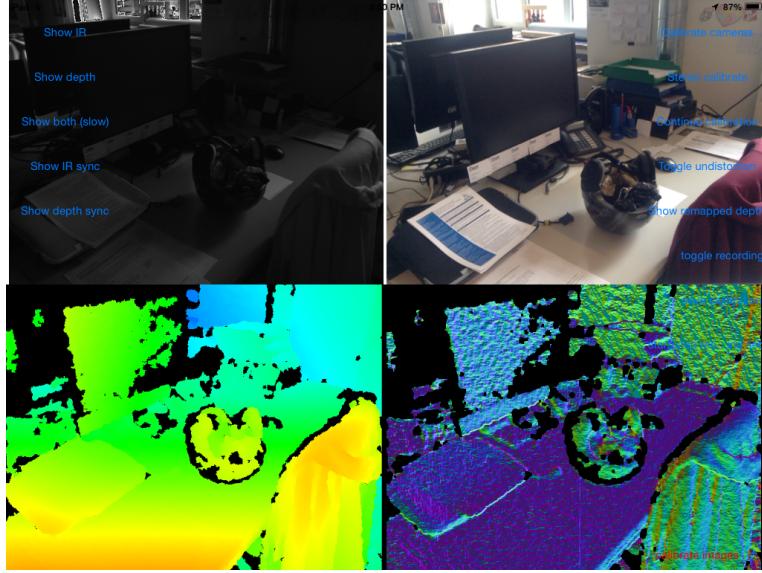


Figure 2.2: Data acquisition screen

2.2.1 Pinhole Camera Model

The pinhole camera model[1] describes the mapping of a 3D point to a 2D image. Equation 2.1 explains the basic relationship between a 3D point and its projection onto the image plane by using perspective transformation.

$$sp = K(RP + t) \quad (2.1)$$

Here $P = (X, Y, Z)^T$ is a 3D point in the world coordinate space and $p = (u, v, 1)^T$ is the projection of P in the camera frame in pixels.

$$K = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \quad (2.2)$$

K is called the camera matrix or a matrix of intrinsic parameters, f_x, f_y are the focal lengths expressed in pixel units and (c_x, c_y) is the principal point which is usually located at the image center. Furthermore we have the rotation matrix

$$R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \quad (2.3)$$

and the translation vector

$$t = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix}. \quad (2.4)$$

We will refer to $[R|t]$ as the matrix of extrinsic parameters. It is used to describe the camera motion around a static scene. In a more general application the extrinsic parameters also describe the rigid motion of an object in front of a stationary camera.

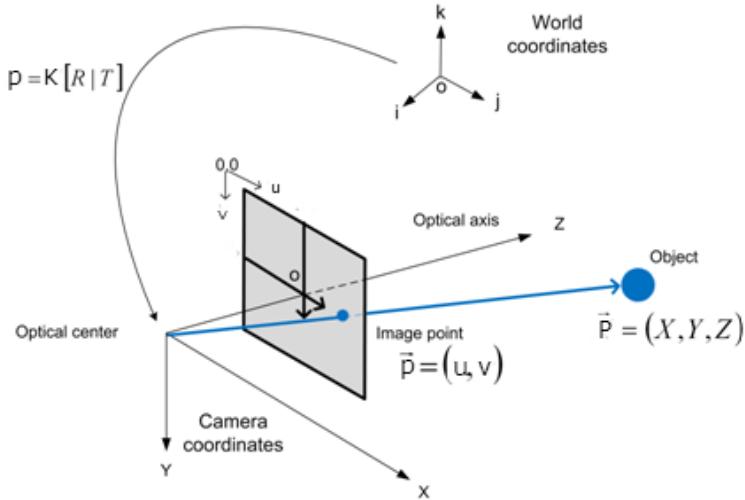


Figure 2.3: Pinhole camera model⁷

Real lenses usually have some radial and tangential distortion. To account for this the distortion coefficients k_1 to k_6 and p_1 and p_2 are introduced and the camera model is extended as follows.

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = R \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + t \quad (2.5)$$

$$x' = \frac{x}{z}, y' = \frac{y}{z} \quad (2.6)$$

$$\begin{aligned} x'' &= x' \frac{1 + k_1 r^2 + k_2 r^4 + k_3 r^6}{1 + k_4 r^2 + k_5 r^4 + k_6 r^6} + 2p_1 x' y' + p_2 (r^2 + 2x'^2) \\ y'' &= y' \frac{1 + k_1 r^2 + k_2 r^4 + k_3 r^6}{1 + k_4 r^2 + k_5 r^4 + k_6 r^6} + p_1 (r^2 + 2y'^2) + 2p_2 x' y' \end{aligned} \quad (2.7)$$

where $r^2 = x'^2 + y'^2$

⁷http://archimede.bibl.ulaval.ca/archimede/fichiers/25229/25229_35.png

$$\begin{aligned} u &= f_x * x'' + c_x \\ v &= f_y * y'' + c_y \end{aligned} \quad (2.8)$$

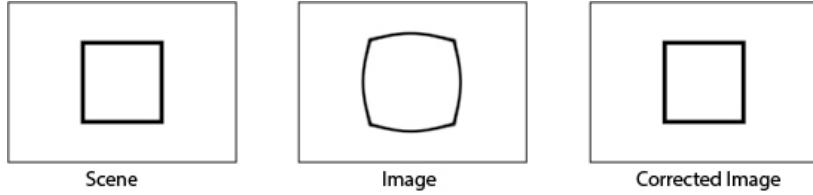


Figure 2.4: Distortion

2.2.2 Single Camera Calibration

We record a set of images of a specific pattern from different viewpoints to do the camera calibration. We choose a chessboard pattern which is a standard choice for calibration because corners can easily be detected. A corner detection algorithm can search for a specified amount of corners, which are the projection of the scene points on the image plane. Because we assume a chessboard we know the location of the corners in a respective coordinate system (which is assumed to be the world coordinate system for the calibration). Accordingly we have several images and thus several equations and we can solve for 2.1 and find 2.2, 2.3 and 2.4.

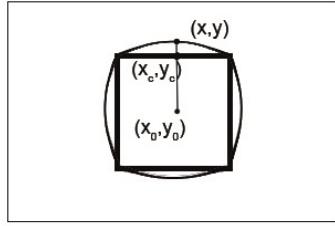


Figure 2.5: Undistortion

Furthermore the calibration process will also find the distortion coefficients by minimizing

$$f(k_1, k_2, \dots) = \sum_i (x'_i - x_{ci})^2 + (y'_i - y_{ci})^2 \quad (2.9)$$

where

$$\begin{aligned} x_c - x_0 &= L(r)(x - x_0) \\ y_c - y_0 &= L(r)(y - y_0) \end{aligned} \quad (2.10)$$

and

$$L(r) = 1 + k_1 r + k_2 r^2 + \dots \quad (2.11)$$

We do this calibration separately for the color camera and the infrared sensor. Typically a good calibration has a RMS of the back-projection error of about 0.1 to 1 pixels. We obtained a RMS of around 0.1 to 0.2 pixels for the color camera and the infrared sensor. We were able to keep the error low by for example fixing the principal point of the image.

We know found the camera matrix and the distortion coefficients for both cameras. In our case this is only an intermediate result, we will continue with the stereo calibration where the single camera calibration results can be used as input parameters.

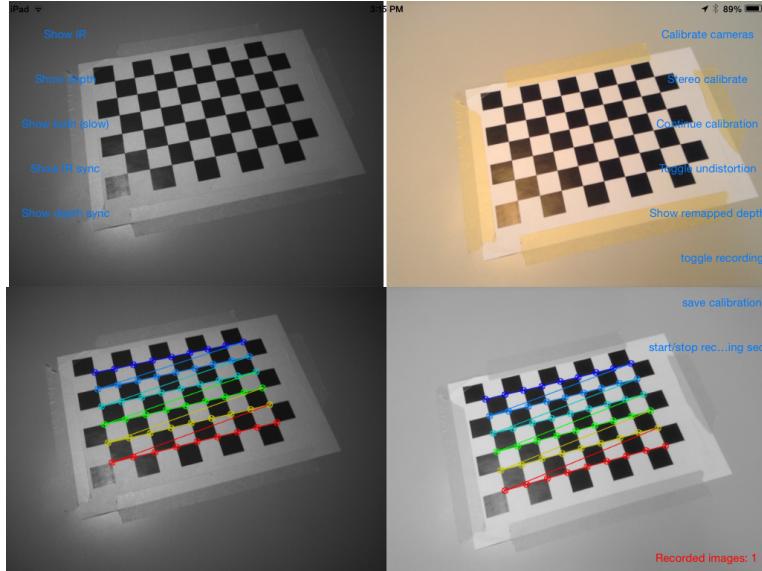


Figure 2.6: IR and RGB images with calibration pattern

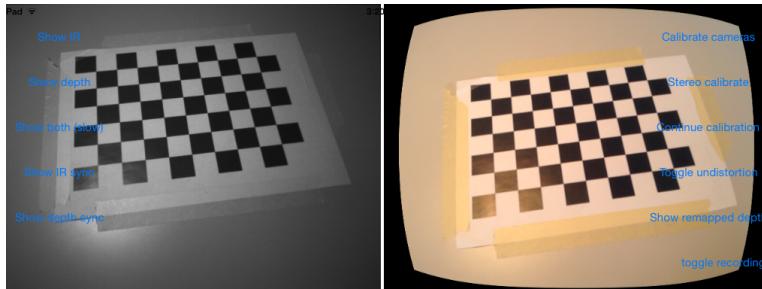


Figure 2.7: Infrared image and undistorted RGB image

2.2.3 Stereo Calibration

In the stereo calibration the extrinsic parameters are estimated. We want to find the rotation and translation between the two camera centers. We will use the precomputed camera matrices and intrinsic parameters from the previous section as well as synchronized RGB-D pairs as input.

Stereo correspondence is usually found using epipolar geometry. Two cameras recording the same scene from different viewpoints span the epipolar plane. It is bounded by the line from the camera center c_0 through the projection of p onto the image plane to p and vice versa for camera center c_1 . The third bounding line is the connection between the two camera centers through the epipoles e_0 and e_1 . Along that, we get a pair of epipolar lines which are the line segments x_0 to e_0 and x_1 to e_1 , extended to infinity.

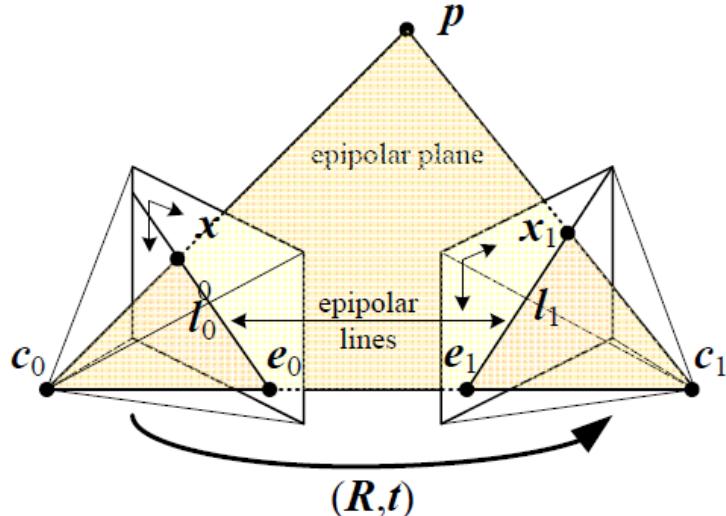


Figure 2.8: Epipolar geometry⁸

The epipolar lines are the intersection of the image planes with the epipolar plane. Then using the epipolar constraint we can define the essential matrix

$$E = [t]_{\times} R \quad (2.12)$$

which can be used to recover the camera motion. The exact procedure is omitted here but using several images it is possible to compute the rotation matrix R and translation vector as described in t [6].

The algorithm uses the calibration points obtained before to construct the epipolar plane. Furthermore we pass the intrinsics computed in the previous

⁸Figure from [6]

single camera calibration as input to reduce the dimension of the parameter space. Additionally flags can be passed to the algorithm, e.g. to fix the intrinsics, to further restrict the solution. The quality of the stereo calibration can be controlled by looking at the RMS like in the single camera calibration. Here the error is usually a little higher but values between 0.5 and 1.5 pixels yielded good results.

2.3 Depth re-projection

The last step in the data acquisition is to use the results obtained in the previous steps and re-project the depth onto the color image. Thus we want to create an RGB-D pair with depth and color information. For storage reasons we will not create a four channel image but use a three channel RGB and a single channel depth image.

From now on, instead of recording infrared/color image pairs, we record depth/color pairs. It must be ensured that the sensor is fixed at the same position relative to the color camera as with which the calibration was performed. The algorithm in A.1 is used alongside with the calibration results to project every pixel in the depth frame to its corresponding location in the color camera image plane.

For every pixel p_{depth} this is done as follows, first we project it into its cameras coordinate system using the respective camera matrix (2.13). Then we apply the rotation and translation obtained during stereo calibration to P_{depth} to transform the point into the color camera coordinate system (2.14). And last we back-project the point P_{rgb} into the color camera frame to obtain its pixel coordinates using the color camera matrix (2.15). Besides the z coordinate of the point in the color camera coordinate system represents the depth of p_{rgb} (2.16).

Mathematically this procedure can be formulated in the following way:

$$P_{depth} = K_{ir}^{-1} \cdot p_{depth} \quad (2.13)$$

$$P_{rgb} = R \cdot P_{depth} + t \quad (2.14)$$

$$p_{rgb} = K_{rgb} \cdot P_{rgb} \quad (2.15)$$

$$\text{depth of } p_{rgb} = P_{rgb}.z \quad (2.16)$$

Because we project the depth values into the color frame we lose some information. E.g. the values which would be projected outside the color frame

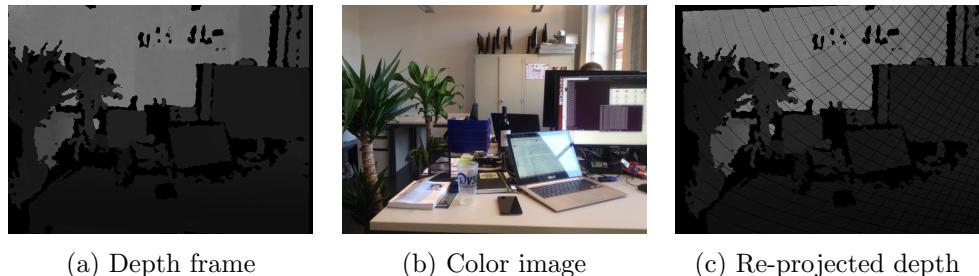


Figure 2.9: Comparison between recorded depth and re-projected depth frame

can't be used anymore, for that reason the back-projected depth frame has black borders. We also see some aliasing effects throughout the entire image which is a sampling problem. This could be solved by a triangulation approach.

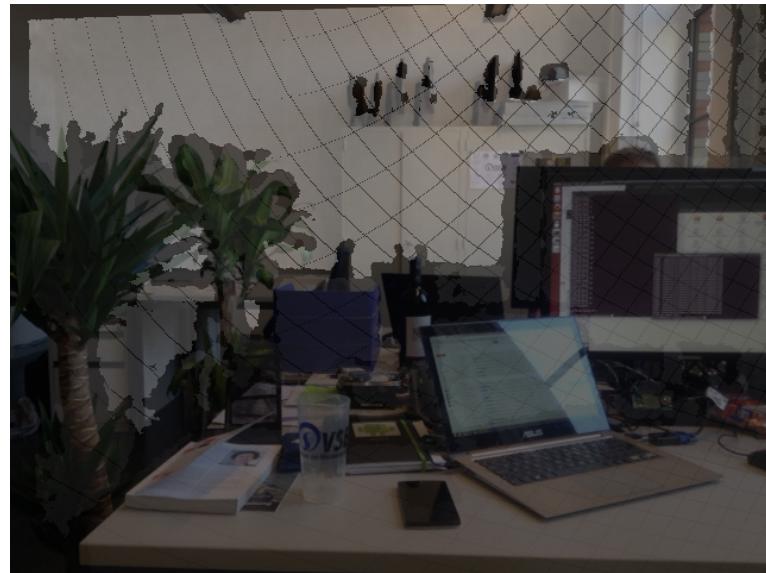


Figure 2.10: RGB-D pair as a 4-channel image

CHAPTER 3

3D Reconstruction

The first step of the 3D reconstruction is to obtain the camera trajectory. For this matter we use a semi-direct motion estimation approach. This pipeline will be slightly extended to output some more results which will be used in the scale estimation step. Due to the fact that the motion estimation is only based on color images we need a scaling step to correct the discrepancy between the absolute metric of the depth map and the relative coordinates of the motion estimation. The final step is to perform the volumetric 3D mapping and texturing which is done by the Fastfusion pipeline.

3.1 Camera Tracking - SVO

For the camera trajectory estimation Semi-Direct Monocular Visual Odometry [2] (SVO) is used. This is a pipeline that implements motion estimation in a semi-direct fashion. It uses two parallel threads, one for estimating the camera motion and the second for mapping.

The motion estimation thread implements the semi-direct approach, which means that it estimates motion directly from intensity values, instead of using robust feature matching as in non-direct approaches. In the first step it initializes the camera pose based on sparse model-based image alignment, it minimizes the photometric error between pixels of different frames projecting to the same 3D location. Then feature-patches are aligned to refine the 2D coordinates of the reprojected points. Finally pose and structure are refined by minimizing the reprojection error introduced in the feature alignment step.

The mapping thread initializes a probabilistic depth-filter for every 2D feature for which the corresponding 3D is to be estimated. When a new keyframe is selected in regions of the image with few 3D to 2D correspondences, new depth-filters are initialized. The filters are initialized with a large uncertainty in depth and then updated at every subsequent frame until its uncertainty becomes small enough, whenever a new 3D point is inserted in the depth map.

As mentioned before the pipeline is extended, it also outputs the $(u, v)^T$

coordinates and the corresponding Z coordinate of every tracked feature. In other words, we obtain the 2D feature's coordinates and the estimated depth. Unchanged is the rest of the pipeline which outputs the camera pose for every frame.

Several challenges needed to be overcome when using the SVO pipeline, because there were several datasets where no camera poses could be estimated. On one side it was discovered that the device movement largely impacted the performance of SVO. Because of compatibility reasons depth maps were saved as XML files on the iPad which strongly influenced the frame rate at which datasets could be recorded. Thus moving the camera and sensor slowly and steady was required. On the other hand good lighting conditions and a complex scene were also crucial to get good results. The preceding point is especially important because the quality of the recorded image was not high (VGA). Hence higher resolution images could improve the results but might need better hardware.

3.2 Scale Estimation

Up to this point we know the camera pose for every frame and the coordinates of the tracked features over all frames. The arising problem is that the estimated 3D coordinates of the features are in a coordinate system with a relative scale. On the other hand the depth values recorded by the structure sensor are saved based on an absolute scale, in millimeters. This means that the relative coordinates must be scaled to match the absolute coordinates.

Thus for a feature at position $(X, Y, Z)^T$ which is projected to $(z_i u_i, z_i v_i, z_i)^T$ in the camera frame, z_i usually doesn't equal the z value of P_{rgb} in 2.14, here denoted as d_i . Therefore we need to introduce a scale factor which corrects the depth difference. We used a simple least squares approach to approximate s :

$$\sum_{i=0}^{pixels \cdot images} (d_i - s \cdot z_i) \cdot z_i = 0 \quad (3.1)$$

Subsequently we scale every camera pose by the factor s . Because the camera pose is computed by

$$c_j = -R_j^T \cdot t_j \quad (3.2)$$

where c_j is the camera center of the j 'th camera and $[R_j | t_j]$ is the camera to world coordinate system transformation matrix. Therefore it is enough to scale the vector t_j by a factor of s to account for the wrong depth.

3.3 Pointcloud

To obtain a visual confirmation that the scale estimation succeeded, the RGB-D pairs were converted to pointclouds. The model from 2.1 is used along with the camera intrinsics to project to the camera coordinate system. Then the rotation and the translation of the camera pose are used as extrinsic parameters and the image is projected to the world coordinate system. Equation 3.3 shows how the pinhole model of 2.1 is used in reverse.

$$P_i = R_{rgb}^T(K_{rgb}^{-1}p_i - t_{rgb}) \quad (3.3)$$

where p_i is the vector $(u_id_i, v_id_i, d_i)^T$. Here u and v are the pixel coordinates in the color camera frame with depth d in the depth frame. Accordingly the color of pixel p_i with depth d is the RGB value at $(u_i, v_i)^T$.



Figure 3.1: Pointclouds visualized in MeshViewer

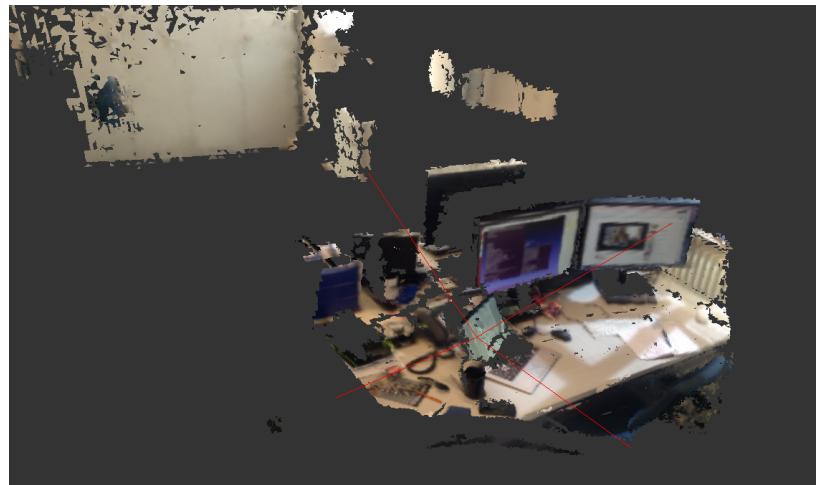
3.4 Dense Reconstruction

The final step of the reconstruction process is the volumetric mapping. Here the Fastfusion [5] pipeline is used because it doesn't need a special sensor, as the Kinect Fusion, thus works with our data. Furthermore it runs on the CPU and doesn't need any special GPU, hence it can be run on any desktop computer or laptop.

Fustfusion uses an octree data structure to perform the volumetric meshing and mapping. Moreover it incorporates two different algorithms, the first for data fusion and the second for incremental meshing, each running in their own thread. The geometry of the reconstruction is represented using a signed distance function (SDF) which is saved in an octree. This choice was made because usually most of the space is either free or unknown. The octree is then able to save different distances at different levels in the tree, thus obtaining several resolutions. In the visualization the highest resolution is used and the lower resolutions serve as fallback. Furthermore every node in the tree contains a small cubic volume called brick which itself contains 8^3 voxels. Lastly the voxel stores the distance, weight and color of a point.

In the fusing process the algorithm iterates over all pixels in an image, computes its 3D point and looks up its corresponding brick in the tree. Upon that the distance value in every concerning brick is updated. In the next step the geometry representation as an SDF needs to be converted to a triangle mesh. This type of storage is better because of a lower memory footprint. Here the Marching Cubes algorithm is applied and mesh cells are created which are associated with branches in the tree.

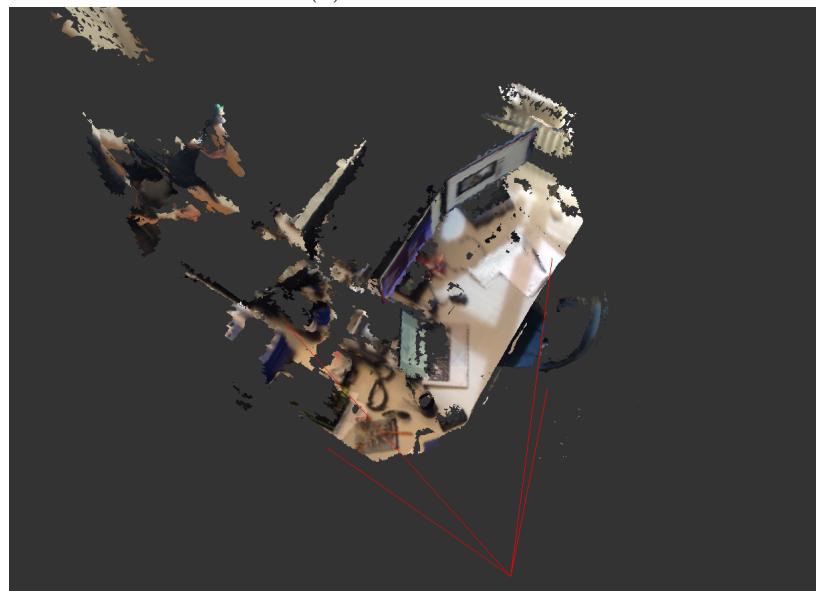
The challenges which were to overcome when using Fastfusion were mostly about the data format. For example the image scale needs to be specified, since the default scale is 5000 (1 millimeter is multiplied by 5000), we chose to also use this format, thus every depth frame needed to be converted accordingly. Moreover the pipeline takes the camera center as in 3.2 and the camera rotation as input. Unfortunately this was not clearly specified, besides the rotation needed to be inverted first which resulted in some confusion.



(a) Desk front view



(b) Desk side view

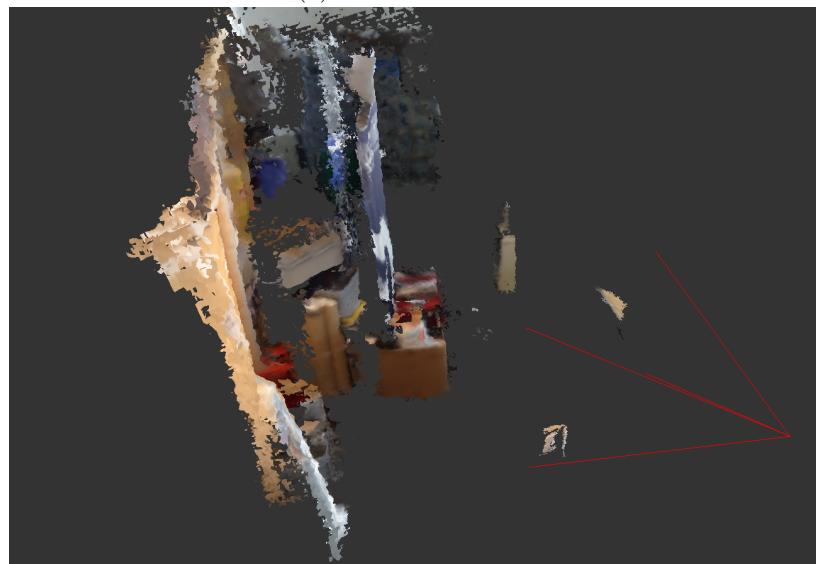


(c) Desk from above

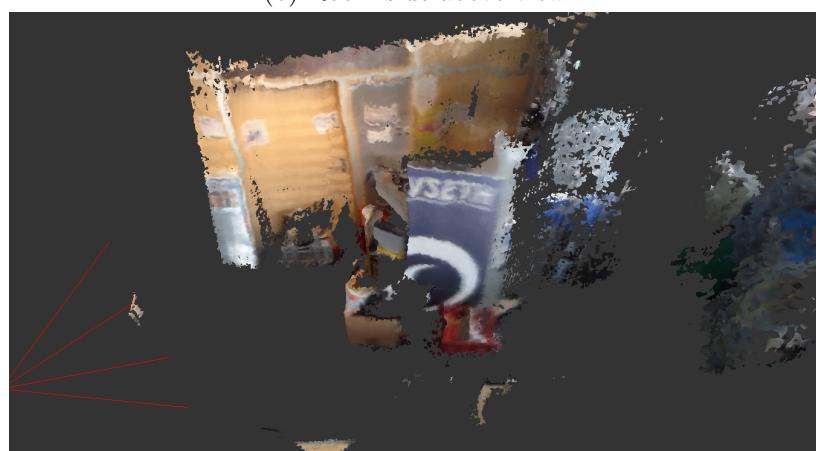
Figure 3.2: Reconstruction example of desk with computers



(a) Room front view



(b) Room side-above view



(c) Room front-above view

Figure 3.3: Reconstruction example of room

CHAPTER 4

App Integration

In this chapter we describe how the depth sensor was integrated into the application described in [7]. The Structure SDK was added to the application such that the depth sensor data can be accessed. The data is streamed live in the background, but depth frames are only saved when a keyframe is selected. Then the reconstruction data consisting of color images, depth frames and the intrinsic and extrinsic parameters for every color / depth pair can be saved on the device and be uploaded to the cloud.

First a visual inertial scale estimation is performed. For this a 3D gyroscope and an accelerometer which come with the device are used, therefore they produce substantial time-dependent and device-specific offsets. Thus the scale needs to be estimated using the current world to camera frame rotation and the current earth-fixed velocity and position are estimated using the inertial sensors. To start the scale estimation, independent motion segments are estimated for visual and inertial based means. After a significant displacement has been determined the two methods are compared and an initial scale is estimated. This procedure is repeated until the scale estimation converges.

Thereupon the inertial position is updated with visual measurements and upon movements of the device the tracking can be achieved through fusion of inertial and visual estimates. The tracking is initialized with two keyframes where ORB features are extracted and matched. After that, relative pose optimization is performed and point matches are triangulated. Upon arrival of a new camera frame, FAST corners are extracted and matched with the previous map. This allows to refine the pose estimate from the inertial sensor and to correct drift.

In our case a new keyframe is added to the map whenever we choose to record one. Upon this we also record a depth frame from the Structure sensor. The keyframe is added to the map and the new observations can be used to match new camera frames. At the same time the depth frame is saved in a buffer for further offline usage. Finally the recorded color images and depth frames are saved to the device combined with the respective pose estimates. In the upcoming offline step we need to perform an additional scale estimation to

account for the discrepancy between the pose estimated using visual and inertial measures and the absolute pose. This step is necessary because the depth frame has the depth values stored in millimeters. Finally we use the same step as in 3.4 to do the 3D reconstruction.

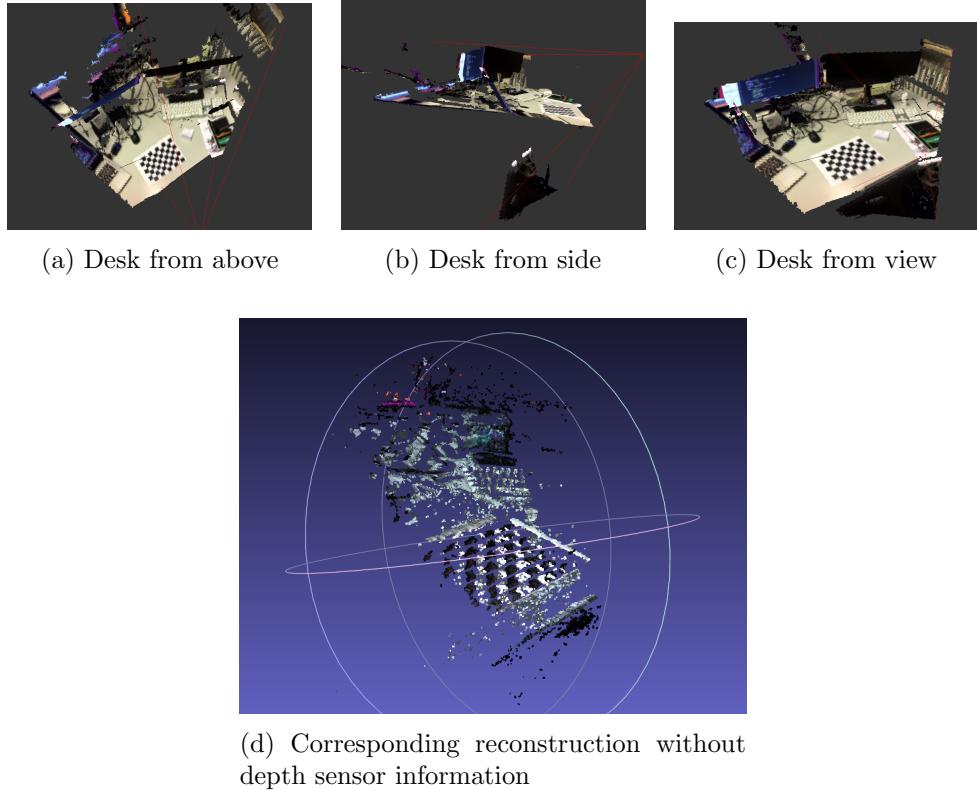


Figure 4.1: Example reconstruction with good results

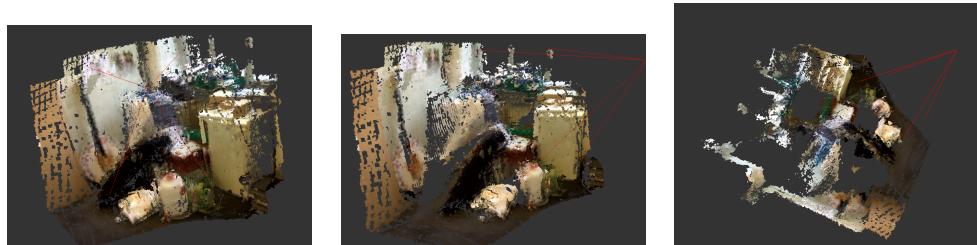


Figure 4.2: Example reconstruction with bad results

CHAPTER 5

Conclusion

The original idea of this thesis was to obtain data from the Structure sensor on a mobile device, calibrate the sensor and the camera to be able to combine the two inputs and finally perform a 3D reconstruction. Additional ideas were to add model fitting or some body pose estimation and extract the skeleton. And there was the possibility to integrate the depth acquisition process into the application of [7].

Whereas we were not able to acquire data on an Android device this worked smoothly on iOS devices due to the available development kit. The main challenge was to extract the recorded data from the device. Especially for the depth frames this was harder than expected because iOS has a different data security model than Android and doesn't allow export of every data format as easily. The resulting workflow of how to transfer the recorded data from the mobile device to the desktop device showed to be easy but could be improved by using a cloud storage or by an all on device approach as described in 6.

The calibration step was first done offline and then ported to the mobile device. Fortunately this worked quite well, the on device calibration has the advantage that the application can easily be run on different devices and there is no additional configuration work to be done offline. Furthermore it has the advantage that the sensor can be moved arbitrarily between different recordings and only a quick two step calibration needs to be performed to update the configuration. Next, the depth re-projection proved to be more cumbersome. First there were different challenges with the algorithm itself, like using the closest depth or eliminating the points which projected outside the image plane. Then re-projecting on the device didn't work, at least not real-time, for performance reasons. So this step needed to be performed offline which is unfortunate in particular with regard to the all on device approach.

Regarding the 3D reconstruction pipeline we note that the original intention was too optimistic. This step showed to be the most challenging and time consuming. Originally the dense visual odometry pipeline should have been used because it conforms with the Fastfusion pipeline. Due to the fact that this pipeline

was running on ROS¹ the usability was hindered. Therefore it was decided to use SVO instead, the drawback being that the camera motion estimation solely relies on color images and not on depth information. As a result the camera motion was estimated in some arbitrary coordinate system, whereas the depth coordinates are stored in millimeters. Therefore scaling must be performed. The approach taken compares the estimated depth values with the real depth and computes an averaging scaling factor. This is then applied to the camera center translation, in practice this proofed to be accurate enough and good results could be obtained. The other drawback of using SVO was that it doesn't work as well in less textured environments. Thus care must be taken when choosing scenes.

The additional tasks like model fitting and body pose estimation were left out because time constraints would not have allowed a profound commitment to the matter. Instead the application described in [7] was made available in iOS and the depth data acquisition could be integrated. This integration nicely completed this work on depth sensors on mobile devices and allowed to show the practicality of mobile depth sensors by enhancing the current application.

Conclusively I can wind up that the original objectives could not all be obtained due to unforeseen challenges and maybe too ambitious prospects. On the other hand different approaches have been tried and therefore and in general a lot about the matter has been learned. Besides, especially the integration of the sensor in an existing application encouraged the enthusiasm for further projects in the field.

¹Robot Operating System

CHAPTER 6

Future work

In this chapter we describe possible future work or extensions to this project. Suggestions are based on efforts which didn't succeed, like for example the Android version. Also tasks which were out of scope because of time constraints will be addressed.

At first we wanted to use Android as platform for the Structure sensor. The problem at the time was that not many devices supported USB host mode which was necessary for the sensor to communicate with the mobile device. Now that more devices support the feature it would be worth to work on data acquisition for Android devices again. The other problem with Android was that the integration with OpenNI and the NDK platform proved to be cumbersome, in the future where Occipital is supposed to release additional resources for Android this could also be solved.

The second point concerns the pipeline. For usability reasons it would be great to have it in one application which runs on a mobile device. Therefore the data acquisition and camera calibration could be combined with a mobile visual odometry or a dense visual odometry approach. This could then be extended to also feature mobile reconstruction. In chapter 4 we described how the previous point was practically solved. It must be added that the future work would be to perform the scale estimation live on the device and to incorporate the recorded depth frames into the fusion thread. In a next step, consumer features could also be added like measuring in the reconstruction or export to 3D file formats to be able to print the recorded object.

An extension to the project which we didn't have time for is model fitting. In the present state everything seen by the camera and sensor is recorded and reconstructed. Still there are cases where you actually only want to record a subset of the viewed scene, e.g. if the background is irrelevant and can be discarded. For that reason we propose to integrate model fitting into the application. The first possibility is to perform a 2D model fitting and then only track and reconstruct the matched region. Another more robust approach is to do 3D model fitting. This has the advantage that there are less features required, thus it works better in less textured environments.

Bibliography

- [1] Docs.opencv.org. *Camera Calibration and 3D Reconstruction OpenCV 2.4.11.0 documentation*. 2015. URL: http://docs.opencv.org/modules/calib3d/doc/camera_calibration_and_3d_reconstruction.html.
- [2] Christian Forster, Matia Pizzoli, and Davide Scaramuzza. “SVO: Fast Semi-Direct Monocular Visual Odometry.” In: *IEEE International Conference on Robotics and Automation (ICRA)*. 2014.
- [3] Msdn.microsoft.com. *Kinect Fusion*. 2015. URL: <https://msdn.microsoft.com/en-us/library/dn188670.aspx>.
- [4] Otmar Hilliges David Molyneaux Richard Newcombe Pushmeet Kohli Jamie Shotton Steve Hodges Dustin Freeman Andrew Davison Shahram Izadi David Kim and Andrew Fitzgibbon. “KinectFusion: Real-time 3D Reconstruction and Interaction Using a Moving Depth Camera.” In: *ACM Symposium on User Interface Software and Technology*. 2011.
- [5] F. Steinbruecker, J. Sturm, and D. Cremers. “Volumetric 3D Mapping in Real-Time on a CPU.” In: Hongkong, China, 2014.
- [6] R. Szeliski. *Computer Vision: Algorithms and Applications*. Texts in Computer Science. Springer, 2010. ISBN: 9781848829350.
- [7] P. Tanskanen et al. “Live Metric 3D Reconstruction on Mobile Phones.” In: *Computer Vision (ICCV), 2013 IEEE International Conference on*. 2013, pp. 65–72. DOI: [10.1109/ICCV.2013.15](https://doi.org/10.1109/ICCV.2013.15).
- [8] Digital Joystick: Video Game News & Views. *E3 2009 : Microsoft at E3 Several Metric Tons of Press Releaseapalloza*. 2009. URL: <http://blog.seattlepi.com/digitaljoystick/2009/06/01/e3-2009-microsoft-at-e3-several-metric-tons-of-press-releaseapalloza>.

APPENDIX A

Code

A.1 Depth Re-projection

```
1  for (int row = 0; row < depthRaw.rows; ++row) {
2      for (int col = 0; col < depthRaw.cols; ++col) {
3
4          // get depth for re-projection
5          double depthVal = (double)ptrOrig[depthRaw.cols*row + col];
6
7          // Map depthcam depth to 3D point
8          Mat<double> P3D = Mat<double>(3,1);
9
10         P3D(0) = (col - camera1.at<double>(0,2))
11         * depthVal / camera1.at<double>(0,0);
12         P3D(1) = (row - camera1.at<double>(1,2))
13         * depthVal / camera1.at<double>(1,1);
14         P3D(2) = depthVal;
15
16         // Rotate and translate 3D point
17         Mat<double> P3Dp;
18         P3Dp = (R*P3D) + T;
19
20         // Project 3D point to rgbcam
21         double xrgb = (P3Dp(0) * camera2.at<double>(0,0) / P3Dp(2))
22         + camera2.at<double>(0,2);
23         double yrgb = (P3Dp(1) * camera2.at<double>(1,1) / P3Dp(2))
24         + camera2.at<double>(1,2);
25         double nDepth = P3Dp(2);
26
27         // "Interpolate" pixel coordinates (Nearest Neighbors)
28         int px_rgbcam = cvRound(xrgb);
29         int py_rgbcam = cvRound(yrgb);
30
31         // Handle 3D occlusions
32         uint16_t &depth_rgbcam =
33         ptrMapped[depthRaw.cols * py_rgbcam + px_rgbcam];
34         uint16_t &depth_rgbcamScale =
35         ptrMappedScale[depthRaw.cols * py_rgbcam + px_rgbcam];
36
37         if(px_rgbcam - depthRaw.cols < 0
```

```
38     && py_rgbcam - depthRaw.rows < 0) {  
39         if(depth_rgbcam == 0 || (uint16_t)nDepth < depth_rgbcam) {  
40             depth_rgbcam = 5*(uint16_t)nDepth; //scale by 5 for fastfusion  
41         }  
42     }  
43 }  
44 }  
45 }
```