

Textured Mesh Reconstruction of Indoor Environments Using RGB-D Camera

Collin Boots

A THESIS

in

Robotics

Presented to the Faculties of University of Pennsylvania in Partial
Fulfillment of the Requirements for the Degree of Master of Science in Engineering

2014

Dr. Daniel D. Lee
Supervisor of Thesis

Dr. Camillo J. Taylor
Graduate Group Chairperson

Abstract

Your abstract goes here... ...

Contents

Abstract	ii
List of Figures	vi
1 Introduction	1
1.1 Motivation and Goals	1
1.2 Related Work	2
1.2.1 Worldview Storage Models	2
1.2.2 Real-Time Processing	5
1.2.3 Additional Resources	5
1.3 Thesis Organization	7
2 Parallel Programming Paradigms	8
2.1 Principles of Parallel Programming	8
2.2 Parallel Algorithm Building Blocks	10
2.2.1 Reduction	10
2.2.2 Scan	11
2.3 Stream Compaction	12

2.4	Programming with CUDA	13
2.4.1	CUDA GPU Architecture	14
2.4.2	Optimizing CUDA Code	15
3	Problem and Approach	19
3.1	High Level System Design	20
3.1.1	New RGB-D Image	20
3.1.2	World Tree	22
3.1.3	GPU World Cache	22
3.1.4	Input Processing and SLAM	22
3.1.5	Image Segmentation	23
3.1.6	Meshing	23
3.1.7	Mesh Improvement	24
4	Implementation	26
4.1	RGB-D Framework Library	26
4.1.1	RGBDDevice	27
4.1.2	RGBDFrame	29
4.1.3	Event Listeners	30
4.1.4	FrameLogger	31
4.2	Filtering and Point Cloud Generation	31
4.3	Plane Segmentation	36
4.4	Planar Mesh Generation	36
4.5	Visualization	36

5	Performance Analysis	40
6	Conclusions and Future Work	41

List of Figures

1.1	Kinect Fusion Tracking and Reconstruction Pipeline. Reproduced from[13]	4
2.1	Parallel Reduction. Reproduced from[22]	11
2.2	Down-sweep phase of parallel sum scan. Reproduced from[22]	12
2.3	Example of stream compacting odd valued integers	13
2.4	How stream compaction works (Scan and Scatter)	14
2.5	CUDA Kernel Organization. Reproduced from [24]	17
2.6	CUDA Memory Organization. Reproduced from [38]	18
3.1	Planned pipeline design. Highlighted blocks represent this thesis.	21
4.1	RGB image with depth overlay. Red is closer to the camera, blue is further away. Represents the output of the RGBD Framework	27
4.2	RGB-D framework architecture	28
4.3	Visualization of surface normals, where components of the normal vector are mapped from -xyz to rgb. Example of output from preprocessing stage.	32

4.4	Preprocessing pipeline diagram.	33
4.5	Comparison of the normals estimated using various methods and filters	37
4.6	Random colorization of detected planes. Example of output from segmentation stage.	38
4.7	Wireframe visualization of the textured mesh. Rendered using the intrinsic parameters of the original camera for direct comparison. Example of output from mesh generation stage.	39

Chapter 1

Introduction

1.1 Motivation and Goals

As robots continue to be incorporated into human environments, the need for intelligent and high-speed reasoning about the objects around them increases dramatically. At the simplest level, mobile robots need to create a map of their environment for navigation. At a higher level, some robots need to recognize distinct objects in their environment, track object movement, and have some intuitive sense of object geometry that is easily stored and processed. Even more importantly, robots must be able to efficiently generate adaptable models of their environment, or worldview, from sensor data in real time. Many different methods for representing the world have been proposed and implemented, and they will be discussed below in more detail. Like the human brain, the robot should also be able to perform these low level functions with only minimal intervention from higher cognitive functions. Such technology also has potential uses beyond robotics. Applications may include easily modeling indoor

environments for interior design concepts, generating 3D tours or maps for various buildings, or low resolution rough mapping of archaeological excavations.

This thesis is intended to work towards such a system based on RGB-D cameras, triangle meshes, and the powerful parallel computing capability modern Graphics Processing Units (GPUs) offer. RGB-D cameras like Microsoft’s Kinect provide a great low cost solution for capturing 3D environments. Triangle meshes are efficient to store, simple to manipulate and refine, and very versatile. Meshes have the added benefit of being well suited to GPU hardware (which was originally designed for just that purpose). This thesis lays out a robust, high-speed GPU based pipeline that converts raw RGB and depth frames into a 3D textured mesh representation of large planar surfaces in the field of view.

1.2 Related Work

A great variety of methods have been applied to RGB-D camera data in an effort to construct a coherent worldview. Each has advantages and disadvantages.

1.2.1 Worldview Storage Models

Point Cloud Models The simplest approach to storing RGB-D data is as a raw point cloud. Each point is stored in a self contained data structure containing at least the point’s position and color information. E.g.

$$P_i = \{pos_x, pos_y, pos_z, red, green, blue\}$$

This approach allows a complete record of the raw data to be stored very easily, but the size of the data stored will grow very quickly, and simply storing the data linearly results in very slow queries.

Nearest neighbor approximations have recently become a popular approach for speeding queries on large point clouds[20]. However, achieving reliably fast queries usually requires some form of hierarchical tree structure like K-d trees[23] or octrees[43]. K-d trees are much more adaptable and usually more efficient in terms of memory storage, but octrees have a significant advantage when it comes to incrementally building a point cloud because points can very easily be inserted into the appropriate octree leaf node with no duplication or restructuring. K-d trees are better suited for compressing point clouds offline.

Voxel Space Models Perhaps the most robust and impressive real-time surface reconstruction algorithm to date is Microsoft’s Kinect Fusion[21, 13]. An open source implementation called KinFu is also available[28]. Kinect Fusion uses a bounded 3D voxel space where each voxel stores the distance to the nearest detected surface or empty space (the default). Figure 1.1 shows the workflow of the Kinect Fusion system. The point cloud generated by each RGB-D frame is projected into the voxel space and each point updates nearby voxels’ distance to nearest surface metric. Implicit surfaces can then be rendered by raycasting through the voxel space and detecting the distance sign crossover point. This results in very high resolution and fidelity reconstructions of the implicit surfaces in the environment. The primary disadvantage of this approach is the workspace size is limited by memory and compute resources which scale with

the resolution and dimensions of the voxel space.

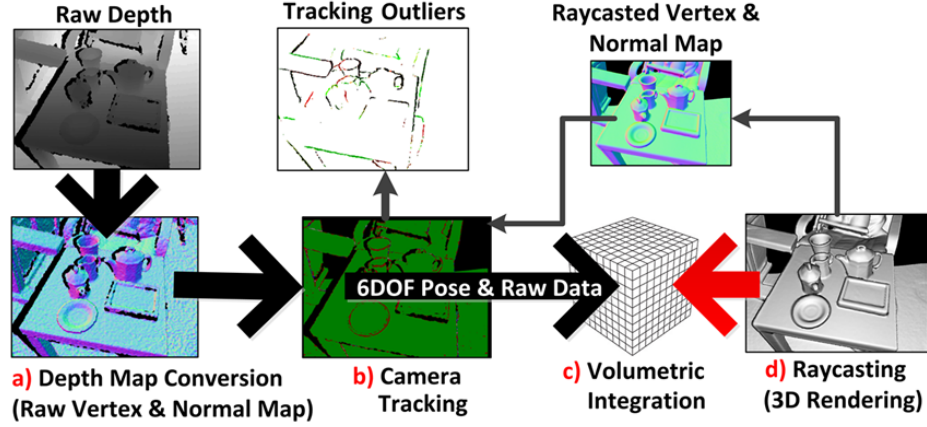


Figure 1.1: Kinect Fusion Tracking and Reconstruction Pipeline. Reproduced from[13]

A spacial extension of this system called Kintinuous was introduced in 2012[42, 41]. Kintinuous uses a mobile voxel space which periodically relocates based on camera motion. Voxels that move outside the space are converted to a triangular mesh using a greedy mesh triangulation procedure described by Marton *et al.*[19].

Others have attempted to use sparse voxel octrees (SVOs) to represent larger unbounded spaces using voxels of varying resolution[16, 31]. This approach is very amenable to human environments like buildings which are dominated by free space. However, all of these approaches are still limited to storing discrete data points and contain no inherent information about coherent objects or regions.

Point Cloud Offline Processing Point clouds can also be processed offline to extract surfaces and objects. These systems are not burdened by the strict requirements of real-time computation so they generally can produce more globally optimal solutions than their online counterparts. These systems generally target unordered

point clouds. Marton *et al.* proposed an incremental approach to triangulation of noisy point clouds[19]. The approach was amenable to introduction of new registered RGB-D frames by only triangulating points that did not correlate with existing mesh models, but the implementation was far too slow for real-time applications. Ma and Others have introduced a system for planar simplification of dense point clouds using a QuadTree based algorithm[18, 17] and texture maps. This thesis will rely heavily on this innovation. Other approaches worked towards implicit surfaces like parameterized smooth surface fits [12] or Poisson surfaces[14, 4].

1.2.2 Real-Time Processing

A crucial part of this thesis is achieving high speed plane segmentation. Many different researchers have created efficient parallel GPU implementations of common segmentation approaches like Markov Random Fields[2], the Potts model[1], parallelized graph-based approaches[40], seeded region growth (SRG)[27], and region growth based on local smoothness constraints[30]. However, this thesis parallelizes an algorithm by Holz *et al.* [10] based on clustering in normal space. This approach is superior for application in this thesis because it can readily detect large planes and solve them globally in a very data parallel manner.

1.2.3 Additional Resources

Although this thesis focuses on detecting and representing only planar elements, the envisioned worldview generating pipeline would need several additional components

and capabilities to be successful. Additionally, many data processing technologies exist that can improve the quality of the input data through more sophisticated filtering.

RGB-D Filtering The Kinect sensor provides its own set of filtering challenges. Khoshelham and Elberink provided a very useful in depth analysis of Kinect accuracy and resolution specifically with indoor mapping applications in mind[15]. Without some amount of preprocessing and filtering, quantization of the depth data and the fact that resolution and accuracy decrease with increasing distance from the sensor would completely prevent this thesis’s pipeline from producing any usable results.

This thesis uses bilateral filtering of the depth image[26, 37]. A simple Gaussian filter is used to smooth local point normals, but other methods exist that could improve results if they could be implemented efficiently in parallel. One such method is adaptively computing filter windows using integral images[11]. Kinect data also is notoriously full of holes from washed out areas, shadows and other noise related effects. Some research makes an effort to fill these holes[5, 33, 44], but this thesis is geared towards progressive improvement of the world model and hopes that any holes will be patched by other viewing angles.

3D SLAM This thesis only deals with processing the current frame, but it was designed with an eye towards creating a closed loop system to integrate data from multiple frames. To accomplish this, a Simultaneous Localization and Mapping (SLAM) system will be needed. SLAM has been practically implemented using a

Kinect using a GPU [35, 34]. Whelan *et al* did fantastic work in comparing multiple SLAM methods combining RGB feature tracking and full point cloud registration algorithms[41]. Other research has evaluated pose graph optimization [6] and point-plane alignment[36] techniques for hand-held 3D sensors like the Kinect.

Mesh Processing and Modification One future direction for this work is to incorporate new information into existing meshes through efficient topology changes and resolution modification. The graphics community offers a wide range of insight into these methods, ranging from tracking surfaces through complex topology evolution[3], modeling deformable solids[32], Mesh subdivision and simplification approaches[29], and texture re-mapping or the effects of image warping[9, 8, 39, 25, 7].

1.3 Thesis Organization

The remainder of this thesis is organized as follows. Chapter 2 will review some basic precepts of parallel algorithm design, along with specific optimization considerations for programming GPUs with CUDA. Chapter 3 will provide an overview of the system design as well as break down the pipeline into sub-modules to be explored in much more detail in Chapter 4. Finally, Chapters 5 and 6 will provide performance analysis, conclusions, and areas of potential improvement for future work.

Chapter 2

Parallel Programming Paradigms

This chapter will provide an introduction to the fundamentals of parallel programming. The remaining chapters will assume a working knowledge of parallel algorithms, GPU architecture, and CUDA optimization techniques.

2.1 Principles of Parallel Programming

As the physical limits of transistor size are being reached, the trend in computing has been away from making processors faster and towards parallel processing. The introduction of General Purpose Graphics Processing Units has made data parallel algorithms a very attractive approach to accelerating computationally intensive tasks. However, migrating from sequential computing to parallel platforms is not always simple or even wise. The following are some principles to keep in mind when considering using GPU acceleration.

Keep data parallel Parallel computing is best suited to performing simple operations on large numbers of independent data elements, such as n-body simulations and per-pixel image processing operations. Aside from these trivial cases, adapting sequential algorithms to a GPU architecture efficiently can be challenging. Not all loops are easily parallelizable.

Memory is a bottleneck Computing with thousands of cores in a GPU makes memory I/O into a major bottleneck. Usually the data transfer to and from the parallel device will take much longer than the actual compute time. Be extremely conscious of memory usage.

Minimize code divergence On most architectures, parallel code runs in batches for increased efficiency. If the code is full of if/then/else statements, the divergent execution paths can dramatically impact performance.

Avoid sequential operations and synchronization Atomic operations that must be performed in order or synchronizations between concurrent threads of execution can be very slow. The more the threads can be independent of each other, the more efficiently the hardware will be able to execute them in parallel.

Be mindful of the hardware When writing performance critical code, understanding the underlying architecture can be a big help. Understanding how the code will actually be distributed and executed across multiple processors will make optimization much easier.

2.2 Parallel Algorithm Building Blocks

Just as there are fundamental building blocks of sequential computing like search, recursion, and iteration, many parallel algorithms make use of primitive operations like reduce, scan, and stream compact.

2.2.1 Reduction

Reduction is any operation on an array of elements that computes a single result. Common examples include sums, average, min, max, and product. In sequential terms, implementing reduction is very simple. Using sum as an example, algorithm 1 shows how this might be implemented.

```
Data: n element array X  
Result: sum of elements in X  
i = 0;  
sum = 0;  
while i  $\leq$  n do  
|   sum += X[i];  
|   i++;  
end
```

Algorithm 1: Sequential Sum

Notice that in the sequential algorithm, the accumulator's value at each iteration depends on the value in the previous iteration. To perform a parallel sum, the associative property of the operator can be exploited to break group the operation into a structured set of binary operations that can be performed in parallel (Figure 2.1).

This parallel reduction method will work for any associative binary operator.

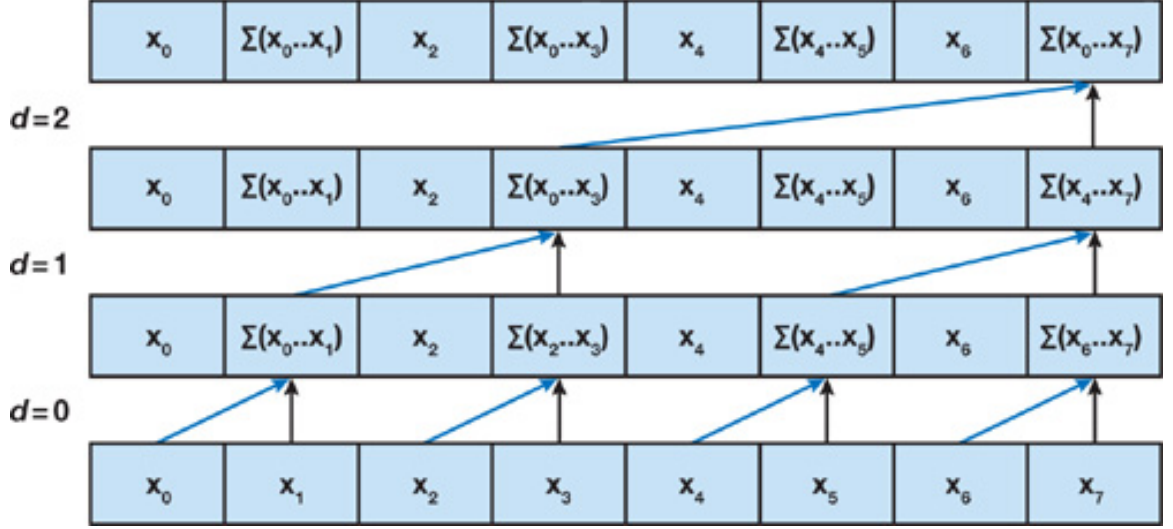


Figure 2.1: Parallel Reduction. Reproduced from[22]

2.2.2 Scan

The scan operation is similar to performing a cumulative sum operation on an array.

There are two forms of scan: exclusive and inclusive. Usually the term scan refers to an exclusive scan.

Given an associative binary operator \oplus with identity value I and an array of n items $[a_0, a_1, \dots, a_{n-1}]$, an exclusive scan returns the array:

$$[I, a_0, (a_0 \oplus a_1), (a_0 \oplus a_1 \oplus a_2), \dots, (a_0 \oplus a_1 \oplus \dots \oplus a_{n-2})]$$

An inclusive scan returns the shifted result:

$$[a_0, (a_0 \oplus a_1), (a_0 \oplus a_1 \oplus a_2), \dots, (a_0 \oplus a_1 \oplus \dots \oplus a_{n-1})]$$

The sequential implementation of a sum scan is trivial and outlined in Algorithm 2.

Data: n element array X
Result: array Y contains exclusive scan of X
 $i = 1$;
 $Y[0] = 0$;
while $i \leq n$ **do**
 $Y[i] = Y[i-1] + X[i-1]$;
 $i++$;
end

Algorithm 2: Sequential Sum

To perform a parallel scan in a work efficient manner, first a reduction is performed on the array as in Figure 2.1. Then, the last element of the array is set to 0, and a series of swaps and sums are performed as shown in Figure 2.2.

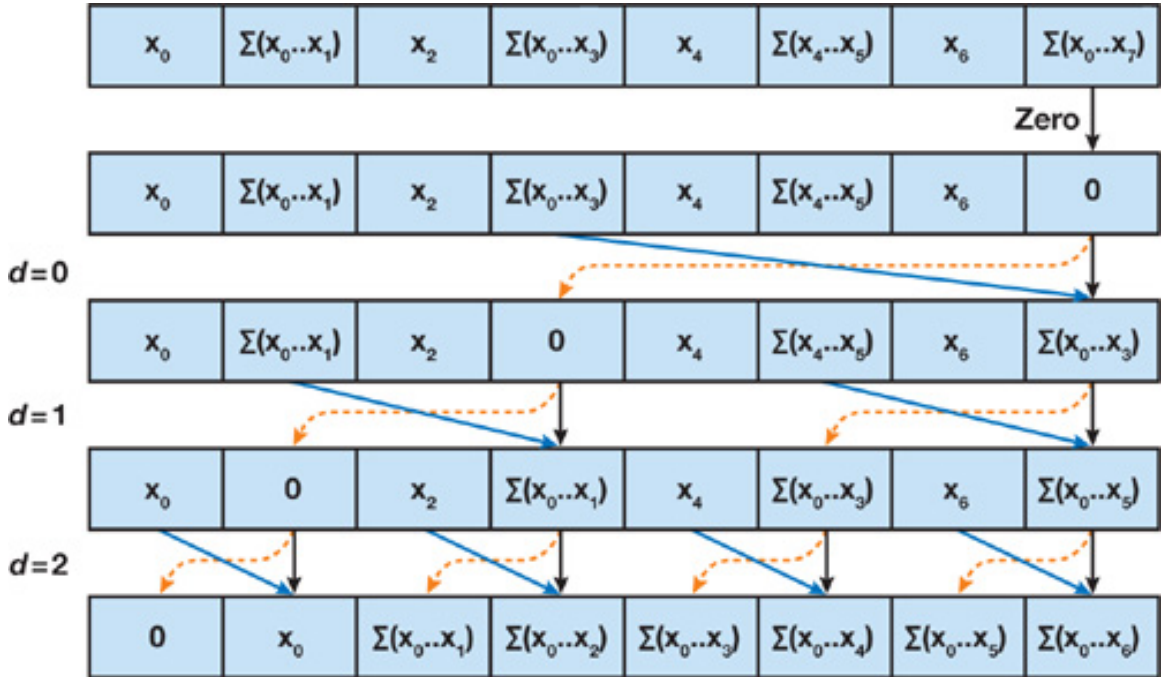


Figure 2.2: Down-sweep phase of parallel sum scan. Reproduced from [22]

2.3 Stream Compaction

In parallel computing, stream compaction is a useful tool for selecting a subset of data and compressing it into a coherent memory block. For example, stream compaction

could be used to transfer all of the odd elements in array A of Figure 2.3 and place them in order compressed at the start of array B.

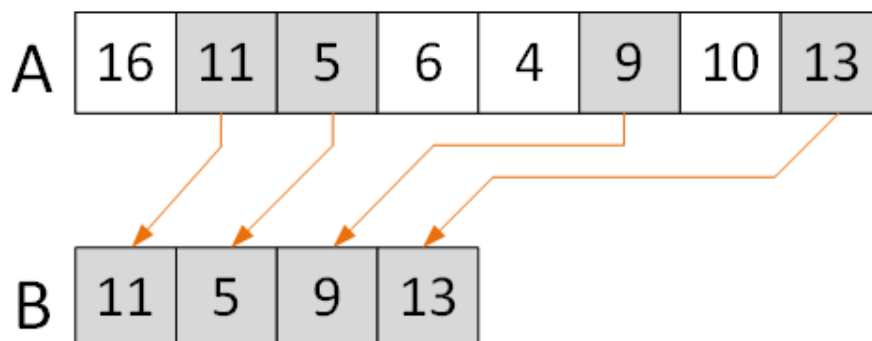


Figure 2.3: Example of stream compacting odd valued integers

Accomplishing this is surprisingly simple. As shown in Figure 2.4, a secondary array of flags is created. This array has a 1 for every odd number in array A and 0 for all others. An exclusive scan is performed on the flag array. Notice that every flagged element now has a unique zero-based index associated with it. The flagged elements are then "scattered" by the computed index. The scatter function is outlined in Algorithm 3.

Data: input array A and scanned flag array F, both of length n
Result: output stream compacted array B
foreach $i=0 \dots n-1$ *in parallel* **do**
 if $isFlagged(A[i])$ **then**
 $B[F[i]] = A[i];$
 end
end

Algorithm 3: Sequential Sum

2.4 Programming with CUDA

CUDA is a C/C++ based programming language created by NVIDIA to expose the general compute functionality of their GPU processors.

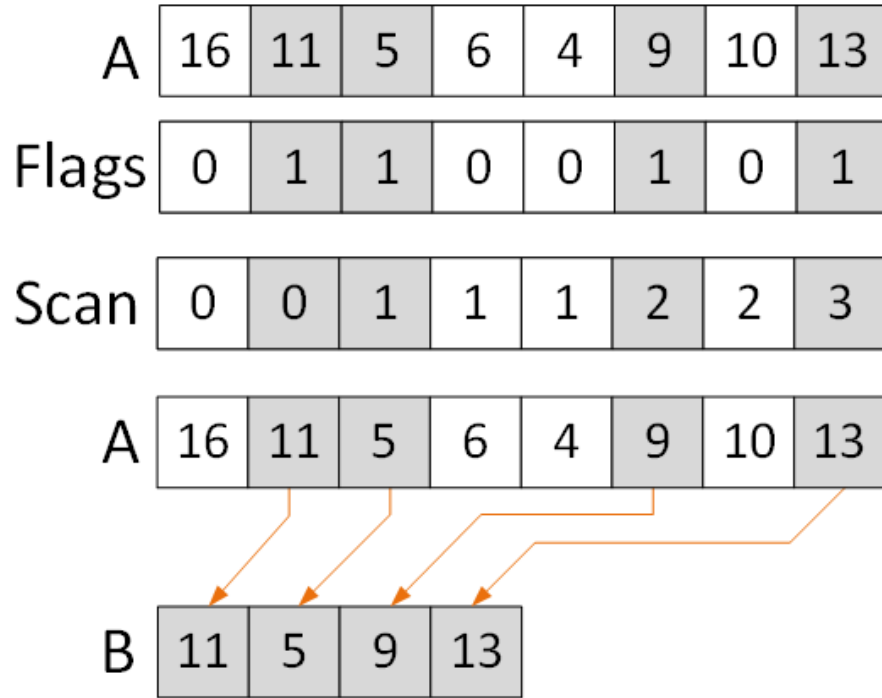


Figure 2.4: How stream compaction works (Scan and Scatter)

2.4.1 CUDA GPU Architecture

CUDA divides the computer into host and device. The host refers to the CPU and all memory it has access to, and the device refers to the GPU and its memory banks. Memory transfers between the host and GPU are started and managed by the host.

A piece of code that runs in parallel on the device is called a kernel. Figure 2.5 shows how kernels are organized. Each individual running process is called a "thread". Threads are organized into "blocks" of up to 1024 threads on modern architectures. Blocks are then grouped into "grids". Within each block, threads are executed in groups of 32 threads called "warps". Every thread in a single warp executes simultaneously on a single processor core. The device memory architecture is depicted in Figure 2.6. Global memory, constant memory, and texture memory are all accessible

by any thread in the kernel. The primary difference between these memories . Every thread also has access to a limited amount of memory which is shared between all threads in a block. Shared memory is much faster than global memory access, and can be used to accelerate kernels that need to share data between threads locally. Finally each thread has access to its own local register memory.

2.4.2 Optimizing CUDA Code

Optimizing CUDA kernels is very application dependent and difficult to provide a universal manual for. However, there are several key underlying GPU architecture features that should be highlighted which can have a dramatic impact on performance.

Global Memory Coalesced Access Global memory is designed in such a way that each half warp can load 64 consecutive bytes in memory as a single parallel read, provided the memory chunk is also 64-byte aligned. This is just enough memory to load an array of 16 consecutively floats at once (one per thread per half warp). To achieve peak memory bandwidth, reads and writes from global DRAM should be grouped into large consecutive chunks of memory. Because of this, it is sometimes extremely beneficial to utilized structure of array (SoA) data storage schemes as opposed to the more typical array of structures model (AoS).

Shared Memory Optimization Shared memory is much faster than global memory, so it can be instrumental in optimizing kernels. For example, instead of performing a reduction or scan in global memory, the entire array can be loaded into shared

memory in a single read operation. The operation would then be performed in shared memory and the results written back in a single coalesced DRAM write. However, shared memory is limited. If too much is used, the number of blocks that can be simultaneously executed on a multiprocessor could go down, preventing the kernel from fully utilizing the GPU hardware. Also, this provides motivation for restricting the size of input data. If all of the data to be processed can be loaded into a single block, shared memory becomes much easier to utilize.

Shared memory is divided into banks, where successive 32-bit words are assigned to successive banks. Shared memory access is done per bank, so if every thread in a warp reads from a different bank, regardless of the ordering, no bank conflicts occur. If two or more threads try to read from the same bank simultaneously, the reads are serialized which reduced the shared memory bandwidth utilization. The exception to this rule is when all threads read from the same bank. This triggers a special broadcast mode that has no bandwidth penalty.

Warp Partitioning and Divergent Code As mentioned above, threads are actually executed in groups of 32 threads called warps. All threads in a warp take the same amount of time to operate. So if only one thread in a warp is actually doing any work, 31 threads will be idle and waste GPU hardware. Minimizing intra-warp code divergence in will go a long way towards fully utilizing hardware. Also, if possible retire entire warps early to free up resources for warps in the execution queue.

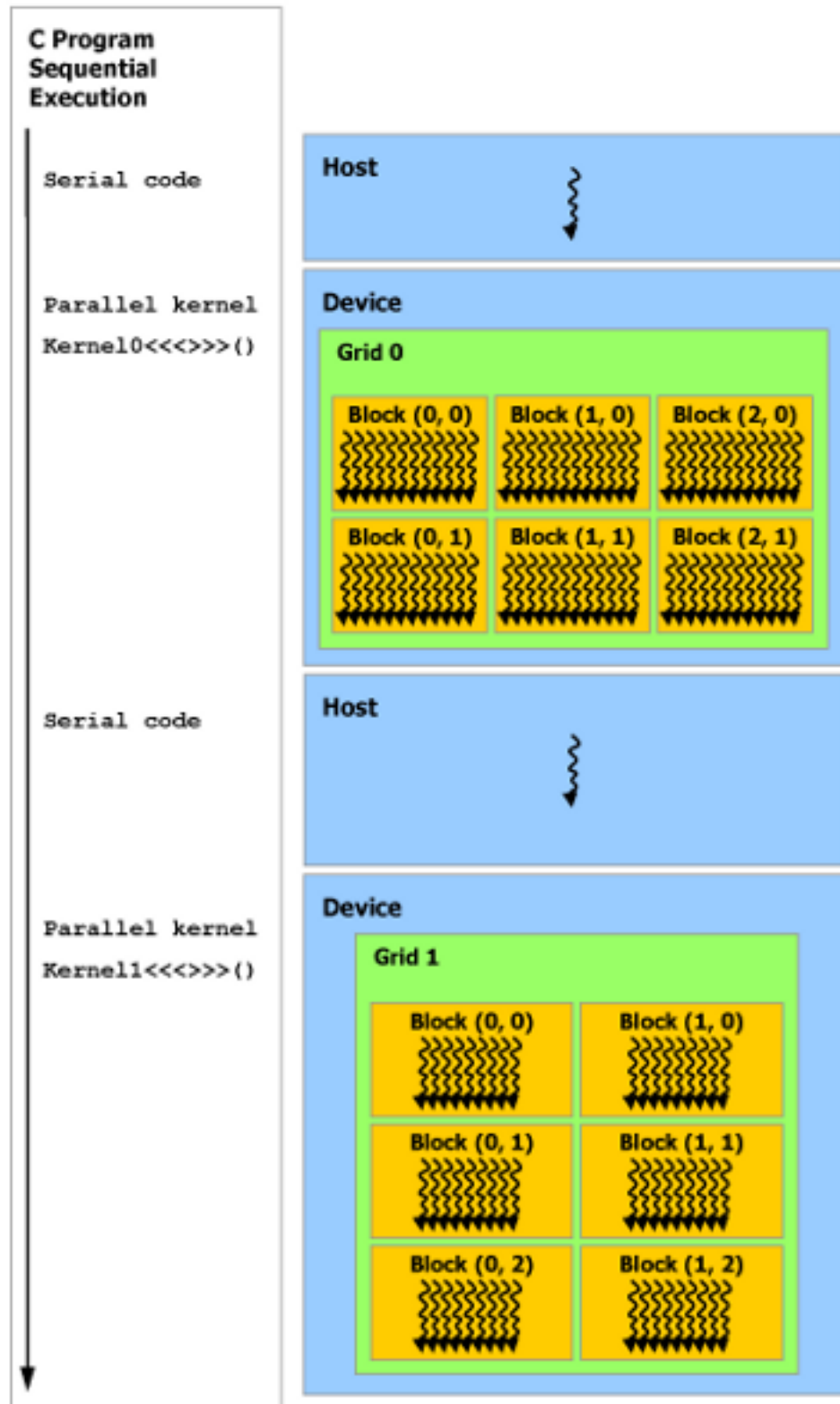


Figure 2.5: CUDA Kernel Organization. Reproduced from [24]

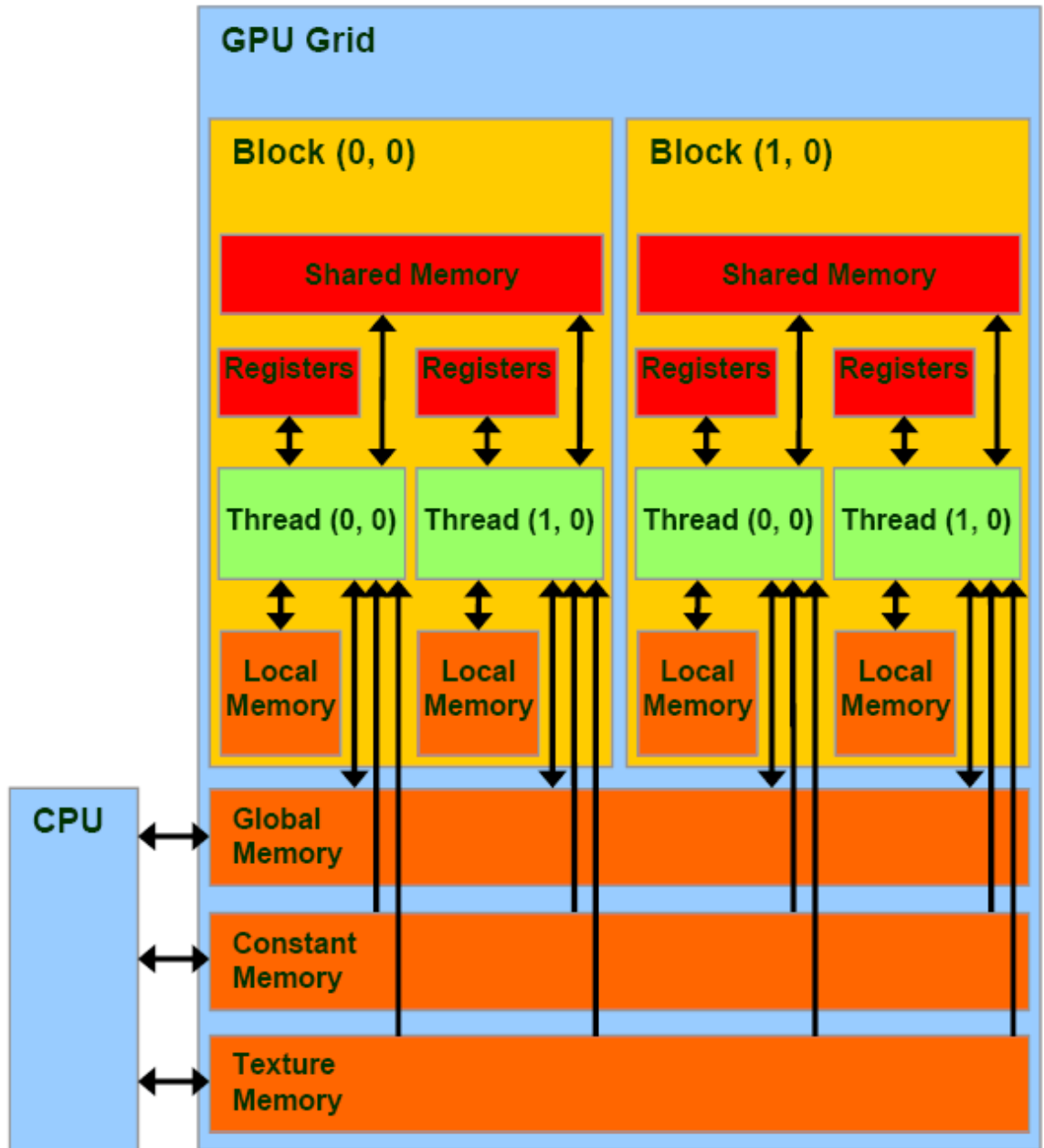


Figure 2.6: CUDA Memory Organization. Reproduced from [38]

Chapter 3

Problem and Approach

Ultimately my goal is to create a software suite that will take RGB-D data from a Kinect or equivalent sensor and add the new information from each successive frame into a mesh-based representation of the world (ideally in real time). The target use case for the system is indoor environments. Human built environments tend to be dominated by planar surfaces, which if segmented out from the image can be processed and stored very efficiently. Since even planar segments with tattered edges require very few triangles to accurately represent them, the hope is that the majority of the scene can be stored with relatively few triangles.

To accomplish all of this in real time, most of the processing will be maintained on the GPU and host-device memory transfers will be minimized. Keeping the bulk of the pipeline on the GPU has many challenges. Existing algorithms have to be heavily modified or re-imagined to run optimally in a data parallel environment. To avoid costly synchronization delays, most of the pipeline must also be designed to accommodate anomalous data without CPU intervention or even acknowledgement.

GPU architecture does not offer a simple way to return results from a kernel, so returning flags for success or failure is difficult to accomplish efficiently.

This thesis presents an implementation of a large portion of the planned pipeline. The pipeline segments out planar segments and converts them to textured meshes using a QuadTree based triangulation scheme.

3.1 High Level System Design

Figure 3.1 shows the planned pipeline at the highest level of organization. The haloed blocks represent the portions of the pipeline that are implemented in this thesis. The other blocks are provided for context, as some of the design features of this thesis would make little sense without the end goal in mind. For the remainder of this chapter I will describe the function of each module in Figure 3.1. The specific implementation of each module will be explained in much greater detail in Chapter 4.

3.1.1 New RGB-D Image

This thesis presents a custom open-source highly modular framework which allows the pipeline to take its RGB-D data from a variety of sources without knowing the implementation details. Currently the framework can only read from a Kinect via OpenNI or from a log generated by the logging tool I developed with the same framework. However, the modular nature of the code allows new data sources to be very easily incorporated.

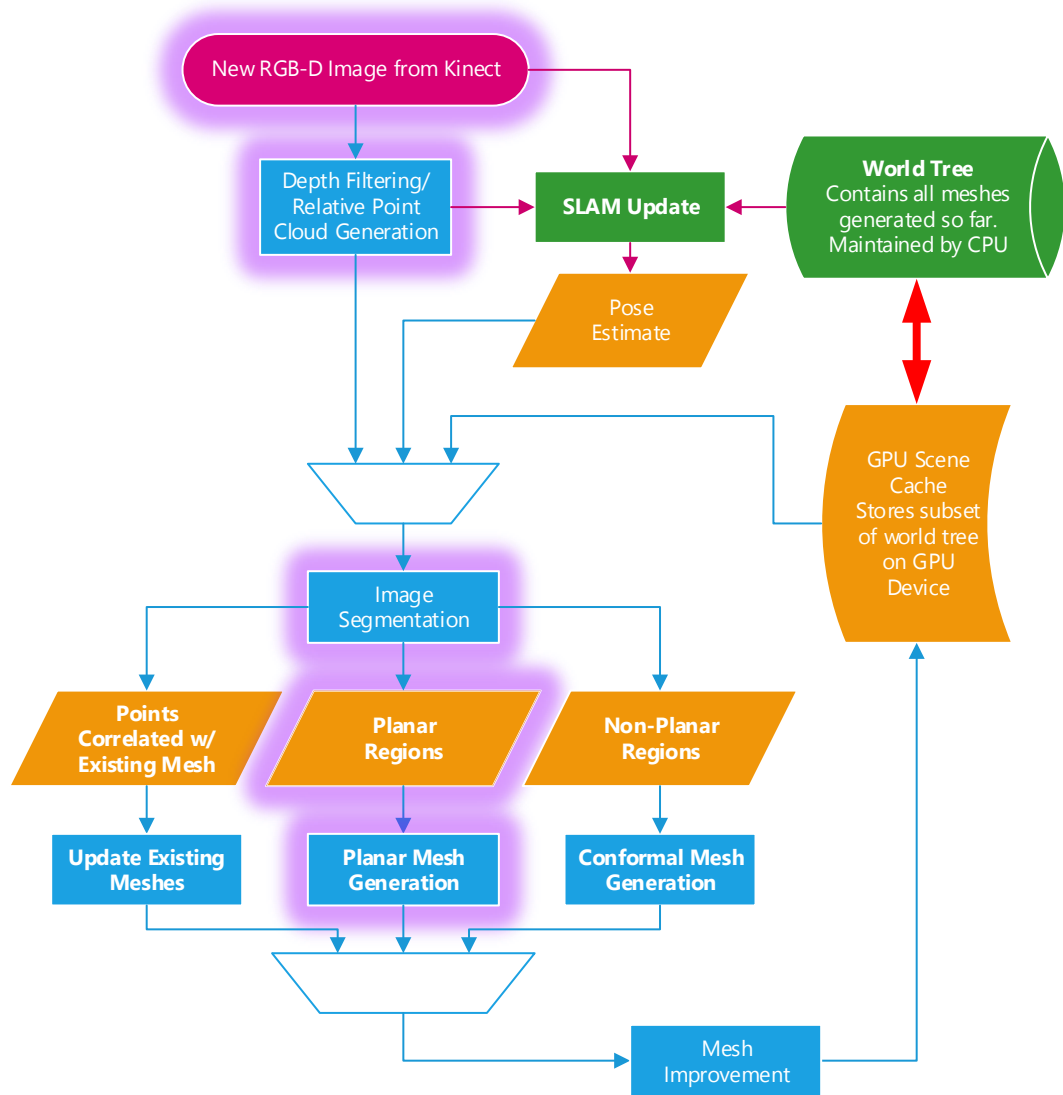


Figure 3.1: Planned pipeline design. Highlighted blocks represent this thesis.

3.1.2 World Tree

The world tree will be maintained in main computer memory by the CPU. It will store the complete mesh representation of the world in a hierarchical tree data structure similar to a scene graph. Unlike a generic scene graph, this structure will be maintained as a strict tree. Each leaf node will contain a single mesh object. Each mesh's component vertices will be stored in an object local reference frame. Objects will be transformed into world space through a series of transformations stored in each node of the tree. This system allows objects to be easily grouped for efficient rendering and simple manipulation.

3.1.3 GPU World Cache

To reduce data redundancy and GPU memory usage, the GPU will only have direct access to a device cached subset of the world tree. The optimal heuristics for maintaining the integrity and efficiency have not been developed yet, but the goal is to maintain a GPU copy of all known meshes in the Kinect's current field of view and those likely to come into view in the next few frames. The cache synchronization with the world tree and subset selection will be handled by the CPU.

3.1.4 Input Processing and SLAM

The raw RGB-D image from the kinect sensor is first passed to the GPU so the data can be filtered and processed in parallel. Once the data is processed, the SLAM algorithm will complete its pose estimation processing, combining the current world

model as stored in the world tree and the 3D features from the newly generated point cloud.

3.1.5 Image Segmentation

Using a combination of information from the processed RGB-D frame, the current camera pose estimate, and the GPU cached world model, each point in the frame’s point cloud will be segmented into three categories: points that correspond to parts of the existing mesh world, uncorrelated planar regions, and uncorrelated non-planar regions. Each category will be processed very differently at the next stage, so generating a robust segmentation will be crucial. Within each category, individual regions will be indexed so they can be easily processed in a region parallel manner for the meshing stage.

For this thesis, only planar regions are detected.

3.1.6 Meshing

Each region generated by the segmentation stage will be processed in parallel and in a very different way.

Correlated Points

Points determined to correspond to existing meshes will be used to update the existing meshes. Points will be used to add new vertices to a mesh as needed for added detail or to increase confidence in existing vertex estimates. This process allows the pipeline to take full advantage of higher resolution/closer capture frames of previously seen

meshes. RGB data will also be used to update model texturing, allowing for texture quality improvement.

Planar Regions

Planar regions are projected to a best-fit true plane and triangulated using a QuadTree-Based (QTB) triangulation algorithm inspired by previous efforts[17, 18]. The method has been demonstrated to efficiently decimate points clouds and dramatically simplifies triangulation and texture generations. In the future, if a planar region is determined to overlap with an existing plane, the existing quad-tree will be expanded to incorporate the new information.

Non-planar Regions

The remaining regions will be triangulated and converted into a new non-planar mesh object using simple polygon based triangulation method. Since conforming texture generation with arbitrary surfaces is a difficult problem to solve in real time with minimal distortion, non-planar meshes will be vertex colored. With a detailed enough triangulation, the visual quality of the of vertex shading is favorably comparable to a textured mesh. The framework would allow for simple addition of textures later if a texture would be a more efficient representation.

3.1.7 Mesh Improvement

Once individual mesh regions have been generated, meshes will be checked for easy optimizations and the potential for mesh merging. The goal of this stage is to reduce

the number of vertices and distinct objects. This step is optional and not required for the pipeline to function, but it can help to simplify the representation and avoid unnecessarily fragmenting the world map.

Chapter 4

Implementation

This chapter will explore in great detail the functionality and implementation of this thesis’s pipeline. As a reminder, the pipeline consists of the highlighted elements in Figure 3.1. Each section of this chapter describes the implementation of one of these functional blocks, as well as a brief exposition on the visualization system. Each section will be accompanied by a figure illustrating the end product of that stage of the pipeline using the same example data.

4.1 RGB-D Framework Library

A crucial component of the pipeline is being able to easily collect RGB-D sensor data from a variety of sources in such a way that the origin of the data is hidden from the remainder of the pipeline. To that end, a highly modular and easily extensible event based framework library was built to seamlessly convert the native data formats and streaming behavior of different sensors. Figure 4.2 provides an overview of

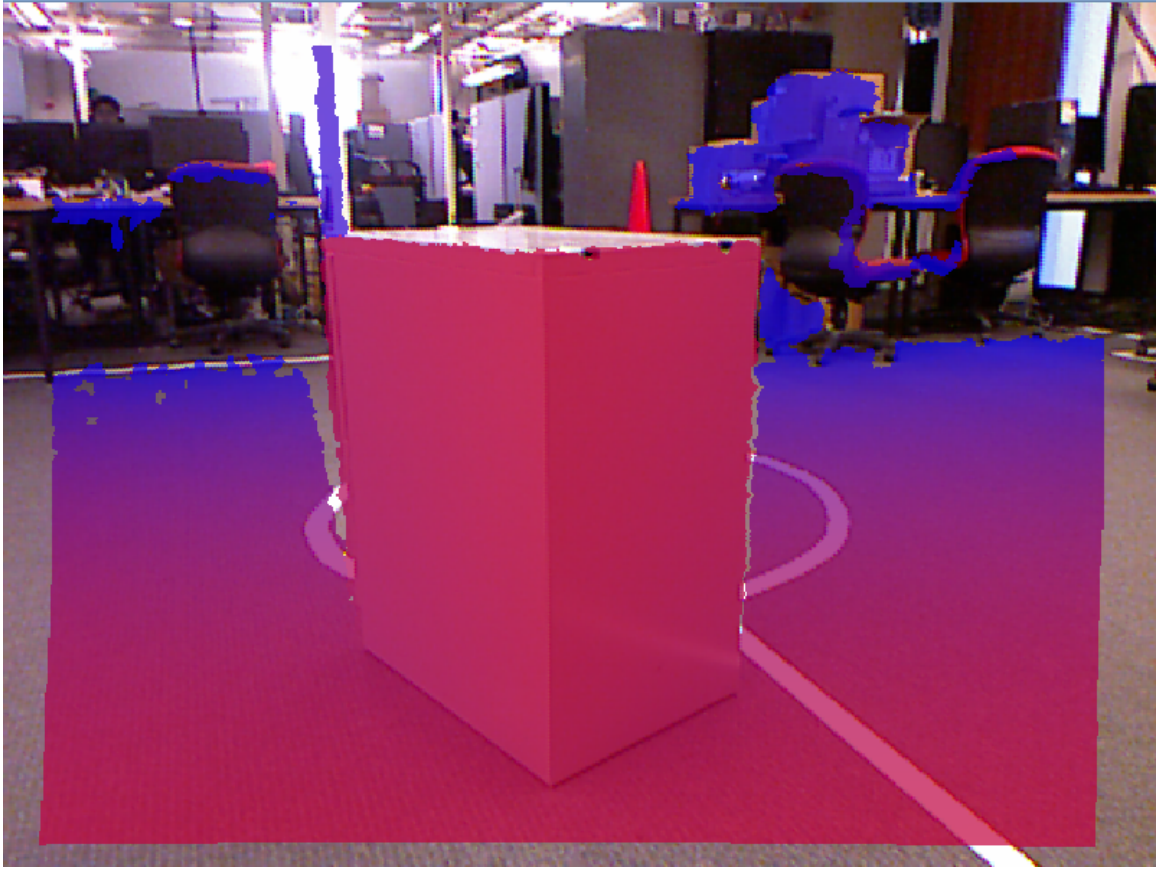


Figure 4.1: RGB image with depth overlay. Red is closer to the camera, blue is further away. Represents the output of the RGBD Framework

the framework organization. The application code deals directly with four primary classes: RGBDDevice, RGBDFrame, Event Listeners, and FrameLogger. The final output of this framework is a color image and a depth image that are registered in such a way that the pixels have a one-to-one correspondence as shown in Figure 4.1.

4.1.1 RGBDDevice

The primary interface between the framework and the rest of the application is the abstract class RGBDDevice. RGBDDevice provides an abstract interface including a data stream management API, access to device properties like camera intrinsics and

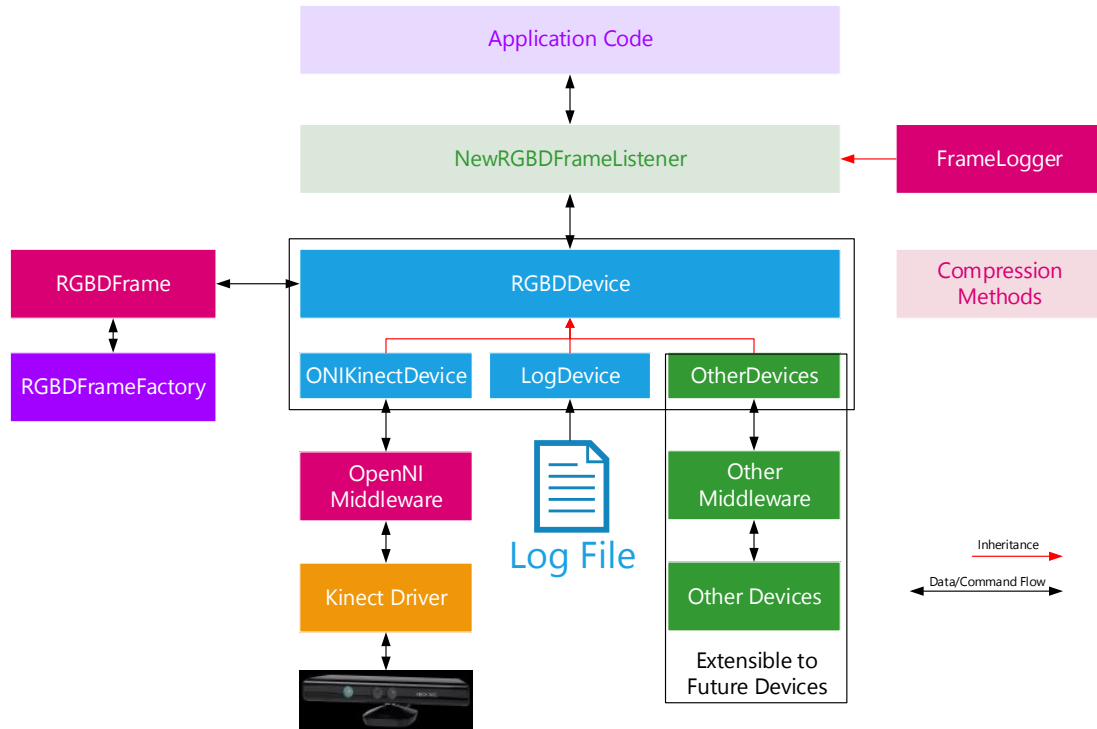


Figure 4.2: RGB-D framework architecture

resolution, and event listener registration API. Interfaces for specific devices like the Kinect can be implemented by creating a subclass of `RGBDDevice` and implementing the abstract methods. The application can then instantiate the desired subclass with device specific initialization parameters and the remainder of the pipeline can be completely agnostic to the underlying nature of the data source, since all `RGBDDevice`s provide the same format `RGBDFrame` through the same event architecture regardless of the native device formatting. The current framework contains two subclasses: `ONIKinectDevice` and `LogDevice`.

ONIKinectDevice `ONIKinectDevice` implements a connection specific to the Microsoft Kinect using OpenNI as middleware. The implementation utilizes OpenNI's

event-based interface to receive data from the sensor. Two streams for the color and depth data are created and registered with listeners internal to ONIKinectDevice. When a new depth or color frame is received from OpenNI, the data is copied and repackaged into an RGBDFrame format and a new thread is launched which passes the RGBDFrame to each registered listener.

LogDevice LogDevice replays a data log created using the FrameLogger class. Because of the way that the FrameLogger records the data, any device that can be implemented as an RGBDDevice can be recorded and played back using a LogDevice. The framework even allows logging data from a LogDevice. This class allows experiments to be performed very easily on prerecorded data without having to alter the behavior of the pipeline at all.

4.1.2 RGBDFrame

The RGBDFrame is a simple data structure that provides a consistent data formatting for users of the framework. Each frame consists of two managed shared pointers (implemented using `boost::shared_array<T>`); one points to the color data and the other points to depth data. Color data is stored in a 24-bit RGB format (1 byte for each component red, green, blue). Depth data is stored in an unsigned 16-bit integer where the least significant bit represents 1mm of resolution (i.e. the value 1024 corresponds to a depth of 1.024m). The frame also has two flags indicating the validity of the depth and color arrays, since not every frame will have both.

RGBDFrames are created using a factory design paradigm. The RGBDFrameFac-

tory generates managed pointers which refer to RGBDFrames of a given resolution. Because of the nature of shared pointers, when the last reference to the frames or its component arrays goes out of scope, the frame will delete itself. Because of this, the frames can easily be handed off to the application with no need for the application to be aware of the finer points of the RGBDFramework's memory management scheme. Since roughly 46MB of frame dedicated memory will be allocated every second by the framework, clean, simple memory management is essential.

4.1.3 Event Listeners

The RGBDDevice provides a suite of event listeners that the application code can register listeners for. Four events can be emitted by the RGBDDevice: DeviceConnected, DeviceDisconnected, DeviceMessage, and NewRGBDFrame.

DeviceConnected This event is triggered when the underlying device is successfully connected to the RGBDDevice. This usually happens upon success of the RGBDDevice::connect method.

DeviceDisconnected This event is triggered when the underlying device is disconnected. This happens either upon calling RGBDDevice::disconnect, or when a fatal communication failure occurs.

DeviceMessage This is a special event that allows the device to pass human readable text descriptors to the application for display in the appropriate place. This is useful for debugging or for registering more unique events for a specific device.

NewRGBDFrame This event actually transfers the latest RGBDFrame from the sensor to the application. It passes a shared_ptr to the all registered listeners.

4.1.4 FrameLogger

The FrameLogger class is actually just a special subclass of the NewRGBDFrameListener class. To record the output of any RGBDDevice to a directory, the FrameLogger just needs to be given an output directory using setOutputDirectory. Once that has been done, all that is required to log the output is to call the FrameLogger's startRecording method, which registers the FrameLogger as a NewRGBDFrameListener with the target RGBDDevice. As the logger receives new frames (which can happen in parallel with the primary pipeline), they are tagged with a counter id and saved to the output directory as raw binary files. Some options are also available for data compression.

4.2 Filtering and Point Cloud Generation

Once the raw RGB-D frame has been received, it is pushed to GPU memory. This is the only downstream memory transfer that happens during this entire thesis. Once the memory transfer is complete, the next step is to convert the RGBDFrame format into a floating point representation that is easier to manipulate and perform other preprocessing steps like depth filtering and normal estimation. The data is also converted to a structure of arrays format to improve memory coherence. Figure 4.4 shows the entire preprocessing system's program flow.

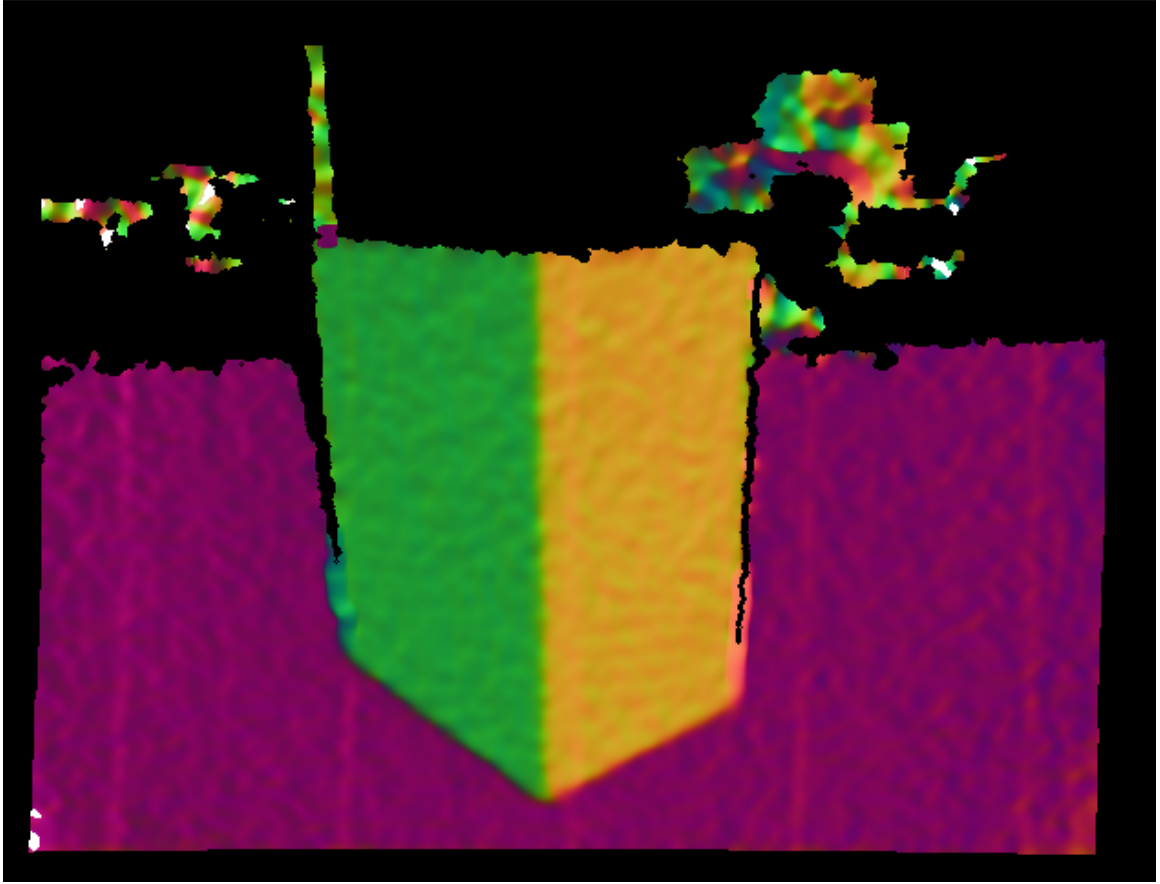


Figure 4.3: Visualization of surface normals, where components of the normal vector are mapped from $-xyz$ to rgb . Example of output from preprocessing stage.

RGB Data Processing The processing for the RGB data is trivial. A kernel converts the 0-255 integer representation of the color data into a 0.0-1.0 floating point representation stored in a structure of arrays (SOA) format for memory coherence. This SOA is then subsampled twice by a factor of two each time, resulting in two additional levels of image resolution. This is useful for later stages in the pipeline that do not require full resolution to function and can therefore be run much faster on the lower resolution image.

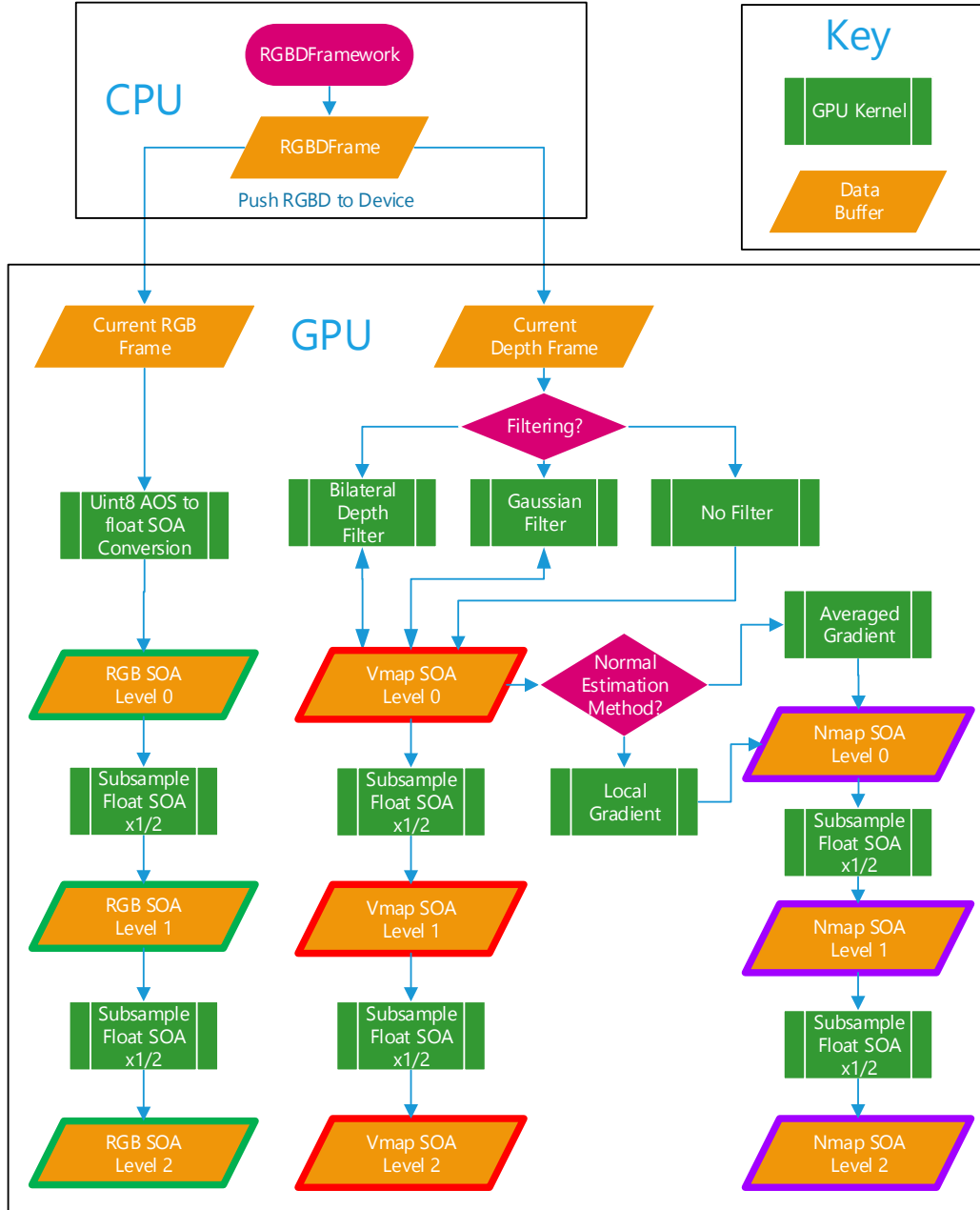


Figure 4.4: Preprocessing pipeline diagram.

Depth Filtering The pipeline includes three options for pre-filtering the depth data: no filter, a separable Gaussian filter, and a separable bilateral filter approximation[26]. The no filter kernel simply performs the inverse camera projection to convert the depth data into an ordered point cloud stored in the vertex map (VMap) buffer. The formula for this projection is based on the camera intrinsic parameters $\{f_x, f_y, c_x, c_y\}$ and the pixel coordinates (u, v) .

$$v_x = (u - c_x) * depth / f_x$$

$$v_y = (v - c_y) * depth / f_y$$

$$v_z = depth$$

Pixels with invalid depth data are stored as NaN to allow easy validation later in the pipeline.

The Gaussian kernel performs the same projection, but also applies a separated implementation of Gaussian filter ($radius = 3pixels, \sigma = 2.0$) to the depth data before calculating the projection. This helps remove much of the depth quantization and sensor noise that makes normal estimation difficult. However, the Gaussian filter does not preserve edges, so crucial discontinuities in the image are blurred together.

To solve this problem, a bilateral filter is used instead. In addition to the spatial Gaussian term that weights pixels by their screen distance from the center pixel, bilateral filters also apply a Gaussian weight to the difference in intensity. In this way, pixels that have a vastly different depth value will have a lower weight and hard

edges can be more easily preserved.

As with the color data, the vertex map is subsampled into a resolution pyramid.

Normal Estimation Surface normals are locally estimated for each valid point. Several methods exist for estimating point cloud normals, but the simplest to apply in an ordered image is a gradient based approach. For each point $\vec{p}(x, y)$ Two vectors \vec{G}_x and \vec{G}_y are created from the neighboring points.

$$\vec{G}_x(x, y) = \vec{p}(x + 1, y) - \vec{p}(x - 1, y)$$

$$\vec{G}_y(x, y) = \vec{p}(x, y + 1) - \vec{p}(x, y - 1)$$

The normalized cross product of these two vectors is then taken to be the point normal.

$$\vec{N}(x, y) = \frac{\vec{G}_x(x, y) \times \vec{G}_y(x, y)}{|\vec{G}_x(x, y) \times \vec{G}_y(x, y)|}$$

Since the sign of this normal is ambiguous, all normals are flipped so that they face the viewpoint. The normal faces the viewpoint when the following condition is met:

$$\vec{N} \cdot (\vec{p} - p_e \vec{y}e) < 0$$

Since in the coordinate frame used at this stage has the eye at the origin looking along the $+z$ axis, this condition can be simplified to:

$$\vec{N} \cdot \vec{p} < 0$$

If the estimated normal violates this condition, the negative of the normal is stored.

While this simple normal estimate works reasonably well, it can be improved by applying a smoothing filter to the gradient images G_x and G_y . A separable Gaussian kernel is applied to each gradient image in turn, and then the cross product is performed. Figure 4.5 compares the various combinations of depth filters and normal estimation methods on the normal estimates. Notice how the Gaussian filter fills in the gaps between the cabinet and the floor, but the bilateral filter shows a very clean break. The final pipeline uses both the bilateral filter and the gradient smoothing techniques.

4.3 Plane Segmentation

4.4 Planar Mesh Generation

4.5 Visualization

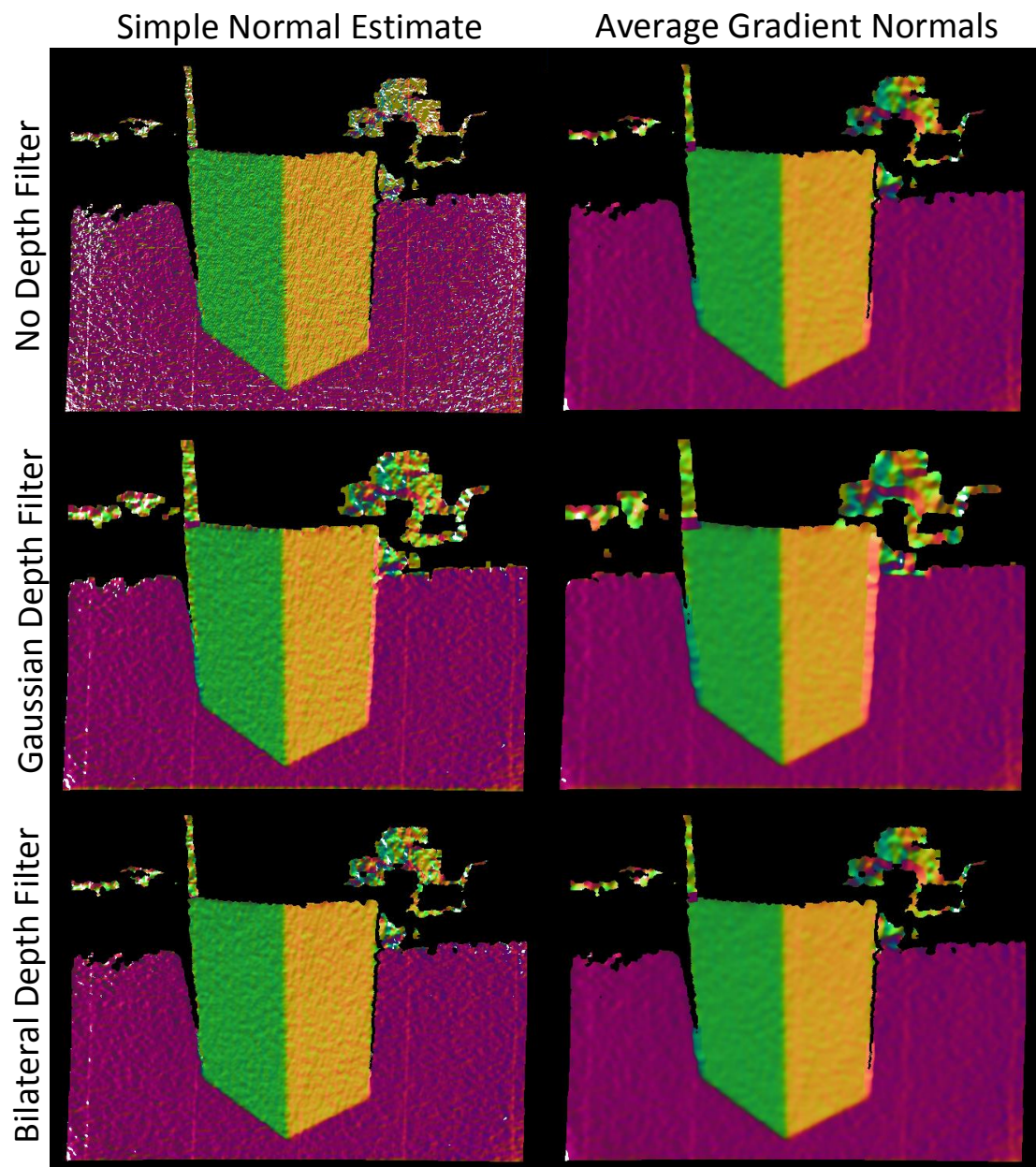


Figure 4.5: Comparison of the normals estimated using various methods and filters

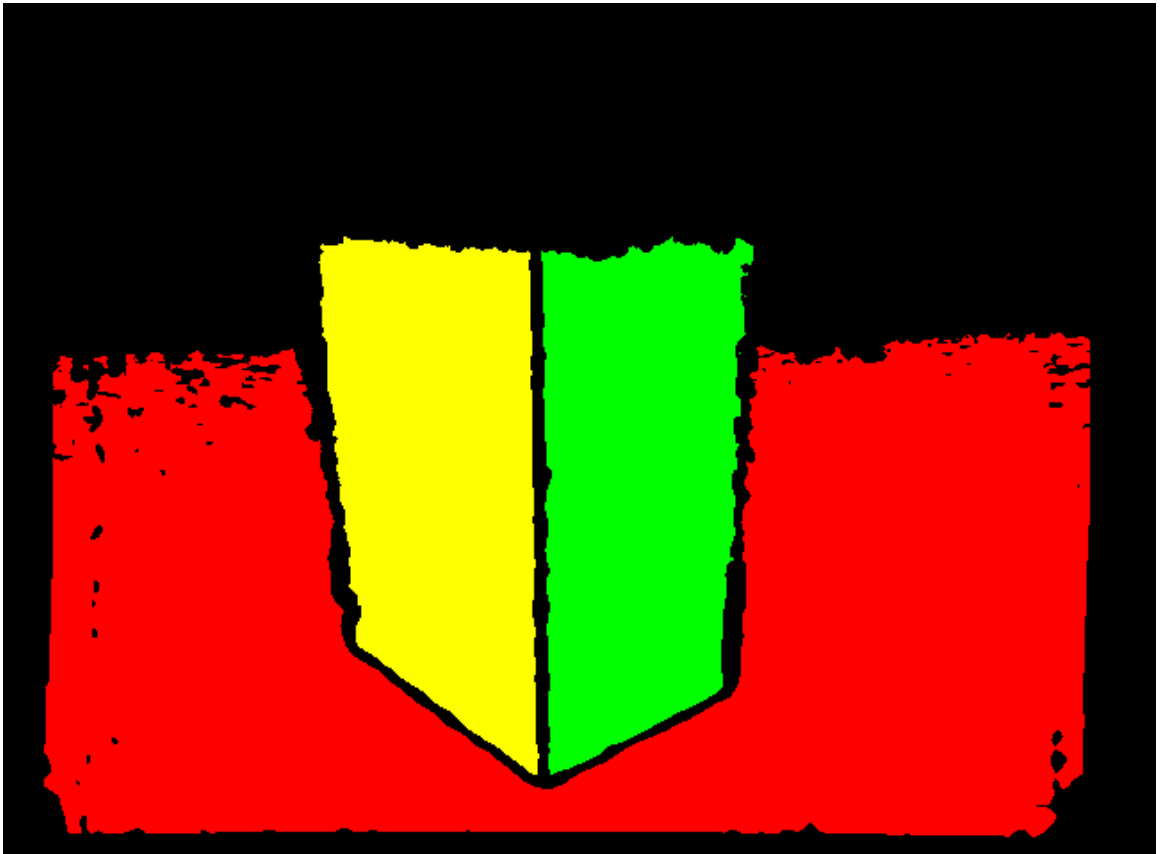


Figure 4.6: Random colorization of detected planes. Example of output from segmentation stage.

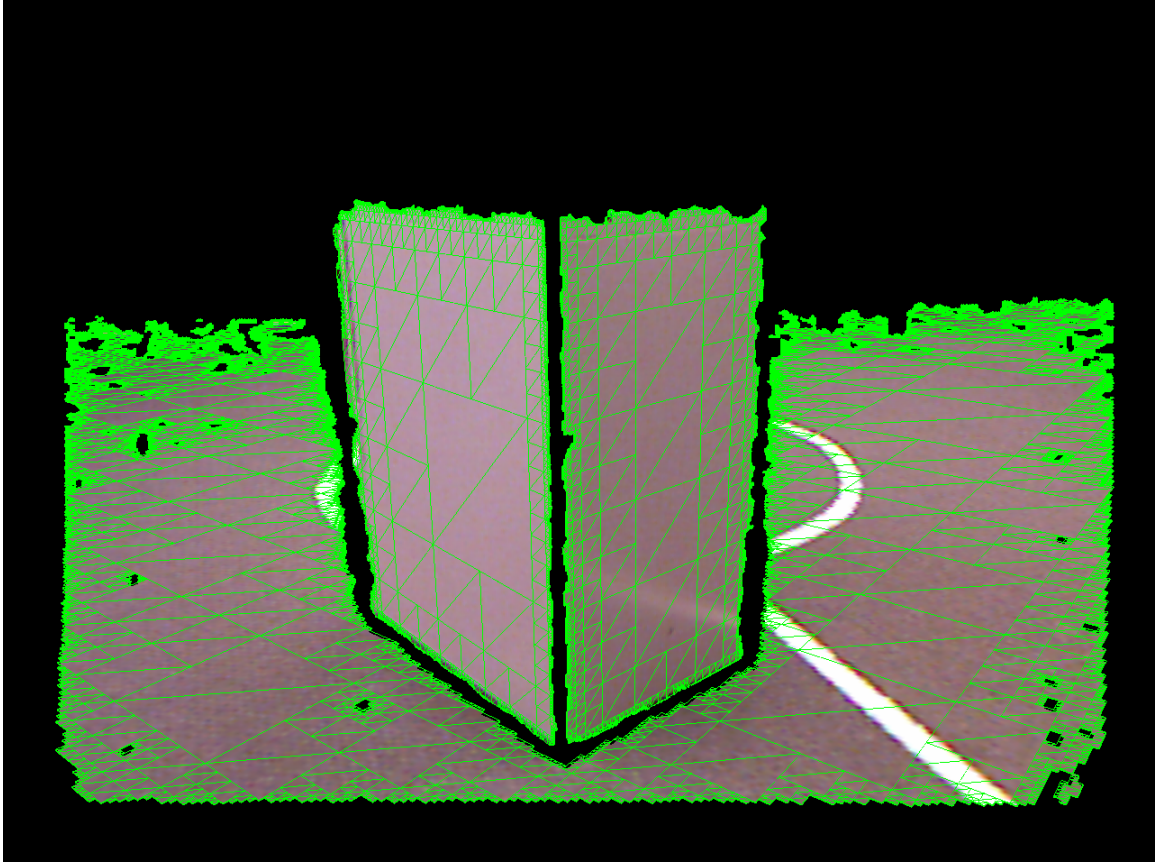


Figure 4.7: Wireframe visualization of the textured mesh. Rendered using the intrinsic parameters of the original camera for direct comparison. Example of output from mesh generation stage.

Chapter 5

Performance Analysis

Chapter 6

Conclusions and Future Work

Bibliography

- [1] Alexey Abramov, Tomas Kulvicius, Florentin Wörgötter, and Babette Dellen. Real-time image segmentation on a gpu. In *Facing the multicore-challenge*, pages 131–142. Springer, 2011.
- [2] Pilar Arques, F Aznar, M Pujol, and Ramón Rizo. Real time image segmentation using an adaptive thresholding approach. In *Current Topics in Artificial Intelligence*, pages 389–398. Springer, 2006.
- [3] Morten Bojsen-Hansen, Hao Li, and Chris Wojtan. Tracking surfaces with evolving topology. *ACM Trans. Graph.*, 31(4):53, 2012.
- [4] Matthew Bolitho, Michael Kazhdan, Randal Burns, and Hugues Hoppe. Parallel poisson surface reconstruction. In *Advances in Visual Computing*, pages 678–689. Springer, 2009.
- [5] Abdul Dakkak and Ammar Husain. Recovering missing depth information from microsofts kinect, 2012.
- [6] Felix Endres et al. An Evaluation of the RGB-D SLAM System. *Robotics and Automation (ICRA), 2012 IEEE International Conference on.*, 2012.

- [7] Yanwen Guo, Hanqiu Sun, Qunsheng Peng, and Zhongding Jiang. Mesh-guided optimized retexturing for image and video. *IEEE Transactions on Visualization and Computer Graphics*, 14(2):426–439, March 2008.
- [8] Paul S. Heckbert. Survey of texture mapping. *IEEE Computer Graphics and Applications*, pages 56–67, November 1986.
- [9] Paul S. Heckbert. Fundamentals of texture mapping and image warping. Master’s thesis, University of California, Berkeley, CA 94720, June 1989.
- [10] Dirk Holz, Stefan Holzer, Radu Bogdan Rusu, and Sven Behnke. Real-time plane segmentation using rgb-d cameras. In *RoboCup 2011: Robot Soccer World Cup XV*, pages 306–317. Springer, 2012.
- [11] Stefan Holzer, Radu Bogdan Rusu, M Dixon, Suat Gedikli, and Nassir Navab. Adaptive neighborhood selection for real-time surface normal estimation from organized point cloud data using integral images. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 2684–2689. IEEE, 2012.
- [12] Kai Hormann. From scattered samples to smooth surfaces. *Proc. of Geometric Modeling and Computer Graphics*, 2003.
- [13] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew 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.
- [14] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, 2006.
- [15] Kourosh Khoshelham and Sander Oude Elberink. Accuracy and resolution of kinect depth data for indoor mapping applications. *Sensors*, 12(2):1437–1454, 2012.
- [16] S. Laine and T. Karras. Efficient sparse voxel octrees. *Visualization and Computer Graphics, IEEE Transactions on*, 17(8):1048–1059, Aug 2011.
- [17] Lingni Ma, R. Favier, Luat Do, E. Bondarev, and P.H.N. de With. Plane segmentation and decimation of point clouds for 3d environment reconstruction. In *Consumer Communications and Networking Conference (CCNC), 2013 IEEE*, pages 43–49, Jan 2013.
- [18] Lingni Ma, Thomas Whelan, Egor Bondarev, Peter HN de With, and John McDonald. Planar simplification and texturing of dense point cloud maps. In *Mobile Robots (ECMR), 2013 European Conference on*, pages 164–171. IEEE, 2013.
- [19] Z.C. Marton, R.B. Rusu, and M. Beetz. On fast surface reconstruction methods for large and noisy point clouds. In *Robotics and Automation, 2009. ICRA '09. IEEE International Conference on*, pages 3218–3223, May 2009.

- [20] Marius Muja and David G Lowe. Fast approximate nearest neighbors with automatic algorithm configuration. In *VISAPP (1)*, pages 331–340, 2009.
- [21] Richard A Newcombe, Andrew J Davison, Shahram Izadi, Pushmeet Kohli, Otmar Hilliges, Jamie Shotton, David Molyneaux, Steve Hodges, David Kim, and Andrew 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.
- [22] Hubert Nguyen. *Gpu gems 3*. Addison-Wesley Professional, 2007.
- [23] A. Nüchter, K. Lingemann, and J. Hertzberg. Cached k-d tree search for icp algorithms. In *Proceedings of the 6th IEEE international Conference on Recent Advances in 3D Digital Imaging and Modeling (3DIM '07)*, pages 419–426, August 2007.
- [24] NVIDIA. Cuda toolkit documentation. <http://docs.nvidia.com/cuda/cuda-c-programming-guide/>.
- [25] Masaaki Oka, Kyoya Tsutsui, Akio Ohba, Yoshitaka Kurauchi, and Takashi Tago. Real-time manipulation of texture-mapped surfaces. In *Computer Graphics (Proceedings of SIGGRAPH 87)*, pages 181–188, July 1987.
- [26] Tuan Q Pham and Lucas J Van Vliet. Separable bilateral filtering for fast video preprocessing. In *Multimedia and Expo, 2005. ICME 2005. IEEE International Conference on*, pages 4–pp. IEEE, 2005.

- [27] Juan C Pichel, David E Singh, and Francisco F Rivera. A parallel framework for image segmentation using region based techniques.
- [28] Michele Pirovano. Kinfu – an open source implementation of kinect fusion. <http://homes.di.unimi.it/pirovano/pdf/3d-scanning-pcl.pdf>, 2012. PhD student in Computer Science at POLIMI.
- [29] E. Puppo and D. Panozzo. Rgb subdivision. *IEEE Transactions on Visualization and Computer Graphics*, 15(2):295–310, March 2009.
- [30] Tahir Rabbani, Frank van den Heuvel, and G Vosselmann. Segmentation of point clouds using smoothness constraint. *International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 36(5):248–253, 2006.
- [31] Julian Ryde and Jason J Corso. Fast voxel maps with counting bloom filters. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 4413–4418. IEEE, 2012.
- [32] Eftychios Sifakis, Tamar Shinar, Geoffrey Irving, and Ronald Fedkiw. Hybrid simulation of deformable solids. In *2007 ACM SIGGRAPH / Eurographics Symposium on Computer Animation*, pages 81–90, August 2007.
- [33] Luciano Spinello and Kai Oliver Arras. People detection in rgb-d data. In *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*, pages 3838–3843. IEEE, 2011.

- [34] F Steinbrucker, Jürgen Sturm, and Daniel Cremers. Real-time visual odometry from dense rgb-d images. In *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, pages 719–722. IEEE, 2011.
- [35] J Sturma, E Bylowb, C Kerla, F Kahlb, and D Cremersa. Dense tracking and mapping with a quadrocopter.
- [36] Yuichi Taguchi, Yong-Dian Jian, Srikumar Ramalingam, and Chen Feng. Point-Plane SLAM for Hand-Held 3D Sensors. Technical report.
- [37] Carlo Tomasi and Roberto Manduchi. Bilateral filtering for gray and color images. In *Computer Vision, 1998. Sixth International Conference on*, pages 839–846. IEEE, 1998.
- [38] Jeremiah van Oosten. Cuda memory model. <http://3dgep.com/?p=2012>, November 2011.
- [39] Lifeng Wang, Sing Bing Kang, R. Szeliski, and Heung-Yeung Shum. Optimal texture map reconstruction from multiple views. In *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, volume 1, pages I–347–I–354 vol.1, 2001.
- [40] Jan Wassenberg, Wolfgang Middelmann, and Peter Sanders. An efficient parallel algorithm for graph-based image segmentation. In *Computer Analysis of Images and Patterns*, pages 1003–1010. Springer, 2009.
- [41] Thomas Whelan, Hordur Johannsson, Michael Kaess, John J. Leonard, and John McDonald. Robust real-time visual odometry for dense rgb-d mapping.

- In *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pages 5724–5731, May 2013.
- [42] Thomas Whelan, Michael Kaess, Maurice Fallon, Hordur Johannsson, John Leonard, and John McDonald. Kintinuous: Spatially extended kinectfusion. 2012.
- [43] K. M. Wurm, A. Hornung, M Bennewitz, C. Stachniss, and W. Burgard. OctoMap: A probabilistic, flexible, and compact 3D map representation for robotic systems. In *Proc. of the ICRA 2010 Workshop on Best Practice in 3D Perception and Modeling for Mobile Manipulation*, USA, May 2010.
- [44] Lu Xia, Chia-Chih Chen, and JK Aggarwal. Human detection using depth information by kinect. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on*, pages 15–22. IEEE, 2011.