

# Enhanced Personal Autostereoscopic Telepresence System using Commodity Depth Cameras

Andrew Maimone\*, Jonathan Bidwell, Kun Peng, Henry Fuchs

*Department of Computer Science, University of North Carolina at Chapel Hill, CB# 3175, Brooks Computer Science Building, 201 S Columbia St., Chapel Hill, NC 27599-3175, United States*

---

## Abstract

This paper describes an enhanced telepresence system that offers fully dynamic, real-time 3D scene capture and continuous-viewpoint, head-tracked stereo 3D display without requiring the user to wear any tracking or viewing apparatus. We present a complete software and hardware framework for implementing the system, which is based on an array of commodity Microsoft Kinect™ color-plus-depth cameras. Contributions include an algorithm for merging data between multiple depth cameras and techniques for automatic color calibration and preserving stereo quality even with low rendering rates. Also presented is a solution to the problem of interference that occurs between Kinect cameras with overlapping views. Emphasis is placed on a fully GPU-accelerated data processing and rendering pipeline that can apply hole filling, smoothing, data merger, surface generation, and color correction at rates of up to 200 million triangles/sec on a single PC and graphics board. Also presented is a Kinect-based markerless tracking system that combines 2D eye recognition with depth information to allow head-tracked stereo views to be rendered for a parallax barrier autostereoscopic display. Enhancements in calibration, filtering, and data merger were made to improve image quality over a previous version of the system.

*Keywords:* teleconferencing, sensor fusion, camera calibration, color calibration, filtering, tracking

---

## 1. Introduction

A long-standing goal [1, 2] of telepresence has been to unite distant workspaces through a shared virtual window, allowing remote collaborators to see into each other's environments as if these were extensions of their own.

In 2002, UNC/UPenn researchers created an early realization of this goal by combining a static 3D model of an office with near-real-time 3D acquisition of a remote user and displayed the result in head-tracked stereo at interactive rates. Since then, several improved 3D capture and display systems have been introduced. In 2004, the MERL 3DTV [3] system offered a glasses and tracker-free capture and display system using an array of 16 cameras and a lenticular autostereo display. However, framerate was low (12 Hz) and the number of viewing zones was limited and repeating. In 2008, the Fraunhofer Institute and the Heinrich-Hertz Institute introduced 3DPresence [4], an improved lenticular-display based system. The system supported multiple views for several participants seated around a table, but like in the MERL system, the number of views was limited and only horizontal parallax was available. In 2009, USC ICT researchers presented a telepresence system [5] that used structured light for 3D acquisition and a volumetric 3D display. The system provided real-time capture, nearly continuous points of view and required no tracking markers or

glasses, but capture and display were limited to a head-size volume. In 2010, Holografika introduced a compelling system [6] consisting of a large array of projectors and cameras offering fully dynamic real-time 3D capture and tracker-less autostereo display. The system, however, featured only a moderate capture rate (10-15 Hz) and did not offer fully continuous points of view – interpolation was performed between a linear array of densely placed 2D cameras and only horizontal parallax was provided. Featuring 27 cameras, 3 PCs, and scores of projectors, it was also a very expensive system to build. In 2011, the FreeCam system [7] demonstrated high quality 3D acquisition using a pair of depth cameras, but capture was limited to users segmented from the background. Also noteworthy are a group of systems [8, 9, 10, 11] with the alternate goal of placing users in a shared virtual space rather than capturing and presenting users within their own physical environments.

In [12], the authors presented a telepresence system that aimed to overcome some of the limitations of previous systems and is the basis for this updated work. The system offered fully dynamic scene capture – presenting a live view of remote users as well as their environments and allowing users to enhance communication by utilizing surrounding objects. Continuous viewpoints were supported, allowing users to look around a remote scene from exactly the perspective corresponding to their head position, rather than from a single or set of fixed vantages. This granted users the ability to see around obstructions and gain more information about the remote scene. Gaze was preserved, allowing participants to make eye contact; re-

---

*Email addresses:* maimone@cs.unc.edu (Andrew Maimone), bidwe@cs.unc.edu (Jonathan Bidwell), pengkun@cs.unc.edu (Kun Peng), fuchs@cs.unc.edu (Henry Fuchs)



Figure 1: Three views of a live capture session.



Figure 2: Two users in 3D scene.

53 search [13] has shown the absence of correct gaze can cause  
 54 a loss of nonverbal communication. Stereo views were pro-  
 55 vided, which have been shown to increase the sense of shared  
 56 presence [13]. Finally, tracking and viewing apparatuses were  
 57 eliminated – 3D glasses obstruct eye contact between partici-  
 58 pants and shutter glasses have been found to be “disruptive” to  
 59 over 90% of users [13]. We believe this system was the first  
 60 to incorporate all of these characteristics – fully dynamic 3D  
 61 scene capture, continuous look-around ability with full paral-  
 62 lax, gaze preservation, and stereo display without the use of any  
 63 encumbrances – but suffered from mediocre image quality. In  
 64 this paper we present a revised system which includes several  
 65 enhancements that result in improved 3D reconstruction.

## 66 2. Background and Contributions

67 Our system is based on the Microsoft Kinect™ sensor, a  
 68 widely available, inexpensive (\$150) device that provides color  
 69 image, infrared image, depth map, and audio capture. Depth  
 70 data is acquired using imperceptible structured light techniques;  
 71 a static dot pattern projected with an IR laser is captured with  
 72 an IR camera and compared to a known pattern [14]. Depth im-  
 73 ages are provided at 640×480 resolution at 30 Hz; color and  
 74 IR images may be captured at this resolution and rate or at

75 1280 × 1024 and approximately 10 Hz. The unit provides a  
 76 58° × 45° field of view and a depth accuracy rated<sup>1</sup> as 1 cm at 1  
 77 m, with a 0.8 m to 3.5 m range<sup>2</sup>.

78 Utilizing several strategically placed and calibrated Kinect  
 79 sensors, an entire room-sized scene can be captured in real-  
 80 time. The scene can be rendered from exactly the remote user’s  
 81 perspective, providing for correct gaze and continuous view-  
 82 points. Eye position tracking is required to provide for contin-  
 83 uous viewpoints; 2D eye detection combined with the Kinect’s  
 84 depth data provides a markerless tracking solution.

85 However, as a device not designed for general purpose 3D  
 86 scene capture or tracking, the Kinect presents some challenges  
 87 for our intended purposes. Since each sensor projects a fixed  
 88 structured light pattern at roughly the same wavelength, inter-  
 89 unit interference is a major problem. The device, as controlled  
 90 with the drivers currently available, provides auto-white bal-  
 91 ance and exposure that cannot be disabled, presenting diffi-  
 92 cultly for seamlessly integrating color-matched data between  
 93 cameras. The capture frame rate, 30 Hz, is suitable for scene  
 94 acquisition but is inadequate for responsive tracking.

95 In [12], we presented solutions to these challenges and in-  
 96 troduced an eye position tracking system based on the Kinect  
 97 that was used with an autostereo display. Specific contributions  
 98 were as follows:

- 99 1. A software solution to the Kinect interference problem  
 100 that provides hole filling and smoothing
- 101 2. A visibility-based algorithm to merge data between cam-  
 102 eras
- 103 3. A visibility-based method for dynamic color matching  
 104 between color-plus-depth cameras
- 105 4. A system for combining 2D eye recognition with depth  
 106 data to provide 3D eye position tracking
- 107 5. A technique for preserving high-quality head-tracked stereo  
 108 viewing on fixed parallax barrier displays even at low  
 109 rendering rates

110 Also presented was a GPU-accelerated software framework  
 111 that implemented hole filling, smoothing, data merging, sur-

<sup>1</sup><http://www.primesense.com/>

<sup>2</sup>We observed that our units return depth readings for surfaces as near as 0.5 m.

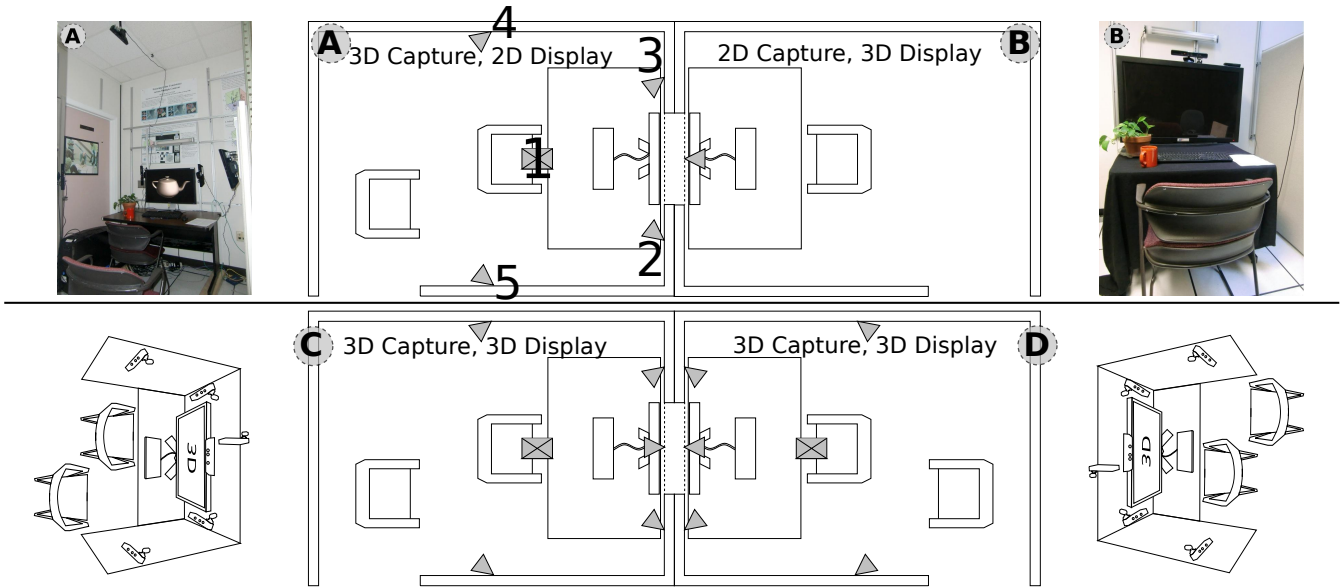


Figure 3: System Layout. Top: Demonstrated proof-of-concept system configuration. Bottom: Proposed ideal configuration.

112 face generation, and color correction at interactive rates for five  
 113 depth cameras on a single PC and graphics board.

114 In this work, we have revised several aspects of the system:

- 115 1. Calibration procedures were enhanced and a new pro-  
 116 cedure was established to correct biases in the Kinect’s  
 117 depth measurements
- 118 2. A smoothing procedure was added to the rendering pipeline  
 119 to suppress depth noise spatially and temporally
- 120 3. The data merger algorithm was updated to be more toler-  
 121 ant of geometric inaccuracies.

122 These enhancements are designed to improve the 3D recon-  
 123 struction quality of the system. In the remainder of this paper,  
 124 we present a complete framework for the system (based on [12])  
 125 with the improvements applied.

### 126 3. System Overview

#### 127 3.1. Physical Layout

128 Figure 3 shows the layout of our system. The two spaces are  
 129 physically separated, but a view of the other space can be seen  
 130 through the display as if the spaces were aligned with a shared  
 131 hole in the wall. The bottom of the figure shows our “ideal”  
 132 configuration – 3D capture and 3D display are supported on  
 133 both sides (spaces C,D).

134 The top of Figure 3 shows the actual configuration used for  
 135 our proof-of-concept system. The system utilizes two office  
 136 cubicles (approximately 1.9 m × 2.4 m). Space A offers 3D  
 137 capture and 2D display of space B, while space B features 2D  
 138 capture and head-tracked 3D display of space A. This config-  
 139 uration allowed us to demonstrate 3D capture, 3D display, and  
 140 eye gaze preservation (for one side) while requiring only the  
 141 single autostereo display that we had available.

#### 142 3.2. Hardware Configuration

143 Both spaces in our proof-of-concept system share a single  
 144 PC with a 4-core Intel Core i7-960 CPU, 6GB of RAM and an  
 145 Nvidia GeForce GTX 580 graphics board. Six Microsoft Kinect  
 146 sensors are connected to the PC. The 2D display side features a  
 147 30 in LCD monitor, while the 3D display side uses a 40 in X3D  
 148 Technologies autostereo display.

149 We avoided networking and audio in our proof-of-concept  
 150 system since both spaces are run from a single PC and are in  
 151 close proximity. We plan to address these omissions in a future  
 152 system.

#### 153 3.3. Software Overview

154 *Operating System and APIs.* Our test system runs on 64-bit  
 155 Linux (Ubuntu 10.10) and uses the OpenNI<sup>3</sup> API along with  
 156 a Kinect driver<sup>4</sup> to communicate with the sensors. OpenGL  
 157 is used for rendering, the OpenGL shader language (GLSL) is  
 158 used for programmable GPU operations, and GLUT is used for  
 159 windowing and user input. The OpenCV<sup>5</sup> computer vision li-  
 160 brary was utilized for camera calibration and tracking.

161 *Data Processing and Rendering Pipeline.* The following ren-  
 162 dering pipeline is used in our system (Figure 4):

- 163 1. When new data is available, read color and depth images  
 164 from Kinect units and upload to GPU.
- 165 2. Smooth and fill holes in depth image.
- 166 3. For each Kinect’s data, form triangle mesh using depth  
 167 data.

<sup>3</sup><http://www.openni.org/>

<sup>4</sup><https://github.com/avin2/SensorKinect>

<sup>5</sup><http://opencv.willowgarage.com/>

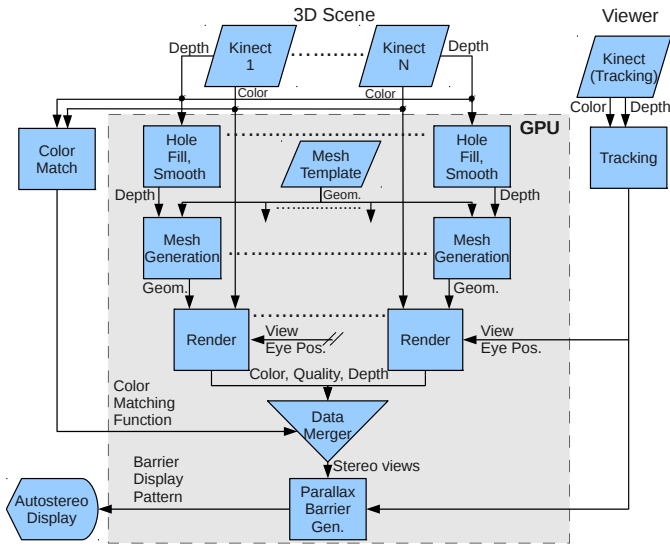


Figure 4: Data Processing and Rendering Pipeline

- 188 4. For each Kinect's data, apply color texture to triangle mesh and estimate quality at each rendered pixel; render from the tracked user's current position, saving color, quality, and depth values.
- 189 5. Merge data for all Kinect units using saved color, quality and depth information.
- 190 6. Repeat steps 3-5 for other eye's view.
- 191 7. Assemble the two stereo viewpoints into pattern required by autostereo display, draw to screen.
- 192 8. While next 3D scene is being generated, periodically re-draw pattern required by autostereo display using the last rendered frame and the new estimated eye position.

180 *GPU Acceleration Overview.* To maximize performance, all real-time graphics-related data processing algorithms are performed on the GPU (using the OpenGL Shader Language) to reduce CPU/GPU memory transfer overhead and take advantage of GPU parallelism. CPU/GPU memory transfers are kept to a minimum: the five Kinects' color and depth images are uploaded to the GPU, but no other major data transfers take place.

187 *System Component Rates.* Since our entire system does not run at the desirable rate of  $\geq 60$  Hz on our current hardware with all graphical enhancements applied, we allow three rates in our system to run asynchronously to allow for interactive rendering and high stereo quality. The first is the *reconstruction rate*, the rate at which new data is incorporated into the rendered scene, which involves uploading new color and depth data to the GPU, hole filling, smoothing, and surface generation. The second rate is the *rendering rate*, the pace at which the scene is rendered from a new viewing perspective. In our system, this also includes our data merger algorithm, which is visibility-based. The final rate is the *parallax barrier pattern generation rate*, the rate at which new patterns are generated for our autostereo display from new estimated eye positions.

201 The independence of the reconstruction rate from the rendering rate helps to keep the system running at interactive rates as more cameras are added to the system; a study by Meehan [15] found a positive correlation between framerate and sense of presence as the former was increased from 10 to 30 fps. The independent parallax barrier pattern generation rate preserves stereo quality during head motion even if the rendering rate decreases. (See Section 5.3.)

209 *Tracking.* We combine 2D eye detection, depth data, and motion tracking to create an unencumbered 3D eye position tracker. Initial eye detection is performed on a color image and eyes are then tracked using pattern matching. Once the 2D eye position is obtained, Kinect depth data is used to transform the position into 3D. A Kalman filter is used to improve accuracy and interpolate the position of the eyes between sensor updates.

## 216 4. Implementation

### 217 4.1. Camera Placement, Calibration, and Error Measurement

218 *Camera Placement.* When placing cameras as shown in the top left of Figure 3, the following factors were considered:

- 220 1. Coverage: for our application, coverage is only necessary for surfaces that can be seen by the remote user.
- 221 2. Redundancy: Redundant coverage allows preservation of surfaces that are occluded from the viewpoint of a single depth camera, but are still visible by the remote user (e.g. an occluded chest behind a raised hand).
- 222 3. Resolution and Accuracy: the resolution available varies with angle and distance (discussion in Section 4.5).
- 223 4. Kinect Depth Range: approximately 0.5 m-3.5 m.
- 224 5. Kinect Depth Interference: discussion in Section 4.2.
- 225 6. Kinect Depth Error: perceived depth error is reduced if camera is near line of sight of user

232 *Camera Calibration.* To calibrate the Kinect sensors, we used the OpenCV camera calibration routines, which are based on Zhang's method [16]. The routines compute camera intrinsic parameters (focal length, center of projection, radial and tangential distortion coefficients) from two or more images of a detected planar pattern, taken from different orientations. Extrinsic parameters (relative positions and orientations) between two cameras were computed using one or more pairs of images of a detected pattern seen from each camera, along with the intrinsic parameters. Since the Kinect driver is able to register the depth image to the color image, only calibration of the color camera is necessary.

244 For our test system, camera intrinsics and extrinsics were computed using detection of a checkerboard target. For extrinsic computation, we calibrated each camera to a master ceiling-mounted camera, whose viewing frustum conveniently overlaps the frustum of each of the other cameras.

249 In addition to the calibration procedures of our previous system [12], we incorporated some of the enhanced calibration procedures of our larger scale 3D telepresence system [17]. To reduce calibration error introduced by the unsynchronized

253 Kinect cameras and motion blur, we placed the calibration tar-  
 254 get at rest before capturing each image. We also used the high  
 255 resolution (1280x1024), low framerate capture mode of the Kinect's  
 256 color camera during calibration and transformed the parameters  
 257 to that of the low resolution (640x480) high framerate mode  
 258 used during system operation. To further reduce error, radial  
 259 distortion in the depth image was also corrected using the dis-  
 260 tortion coefficients measured during intrinsic calibration of the  
 261 color camera, as the two images are registered by the driver.

262 *Depth Bias Correction.* After revising our calibration proce-  
 263 dures, we continued to observe small geometric misalignments  
 264 between Kinects that occurred primary in the cameras' lines of  
 265 sight and hypothesized a bias in the Kinect's depth readings. To  
 266 test this hypothesis, we compared the Kinect's measurements to  
 267 a laser rangefinder (Leica Disto™plus, accuracy  $\pm 1.5$  mm) at  
 268 several distances using the following procedure:

- 269 1. A 2 m sliding rail was placed approximately 1.8 m from a  
 270 wall so that its direction of movement was approximately  
 271 perpendicular to the wall.
- 272 2. The laser rangefinder was mounted on the rail facing the  
 273 wall. The orientations of the rail and rangefinder were  
 274 finely adjusted until the laser spot emitted by the range-  
 275 finder onto the wall remained nearly motionless as the de-  
 276 vice was moved the length of the rail.
- 277 3. A Kinect unit was mounted to the rail so that it was ap-  
 278 proximately perpendicular to the wall, and a checker-  
 279 board target was placed on the wall. The orientation of  
 280 the Kinect was finely adjusted until the calculated extrin-  
 281 sics between the Kinect and checkerboard indicated per-  
 282 pendicularity.
- 283 4. At the closest position on the rail, the z-axis translation  
 284 component of the extrinsics between the Kinect and the  
 285 checkerboard was compared to the corresponding range-  
 286 finder distance to determine their distance offset.
- 287 5. At approximate 0.25 m intervals along the rail, the range-  
 288 finder and Kinect distances were recorded. Kinect dis-  
 289 tances were computed as the mean distance over a  $100 \times$   
 290 100 sampling window near the center of the depth map.

291 The test setup described in steps 1-3 is pictured in Figure 5.

292  
 293 Using the aforementioned procedure on several units, we ob-  
 294 served that real-world distances returned by the utilized Kinect  
 295 driver were biased and tended to overestimate depth linearly  
 296 with distance. We used a least-squares fitting of our measure-  
 297 ments to build a linear depth correction function for each Kinect.  
 298 Table 1 shows the original and corrected measurements for one  
 299 unit.

300 *Depth Error Measurement.* Since our mesh generation and data  
 301 merger techniques rely on knowledge of the relationship be-  
 302 tween the distance of a surface to a Kinect depth camera and  
 303 measurement error, it is beneficial to characterize this relation-  
 304 ship. We expect the depth resolution to fall off quadratically  
 305 with distance.

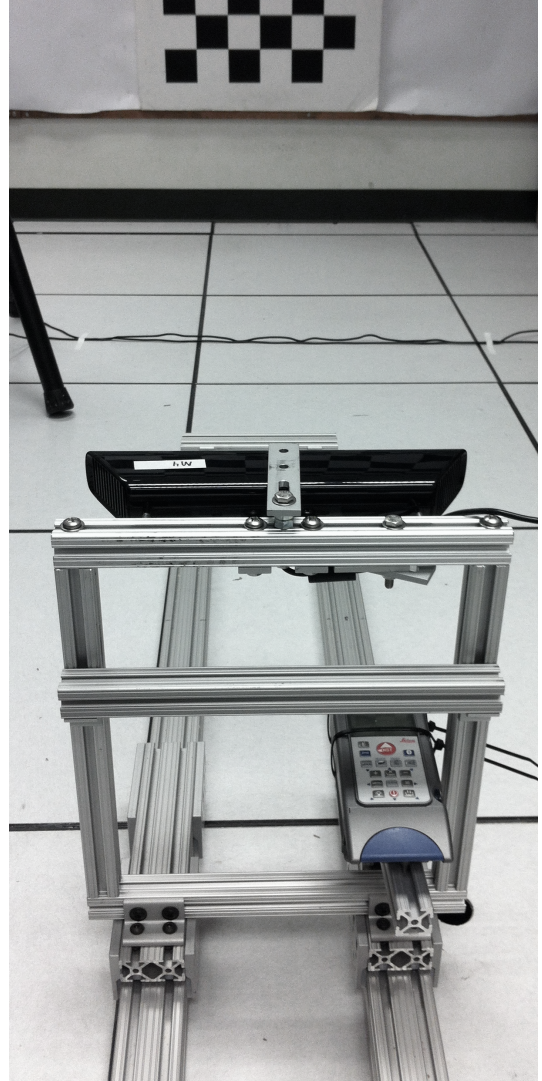


Figure 5: Test setup used to measure depth biases.

Table 1: Typical Depth Bias Correction (all values in mm)

| Kinect | Rangefinder | Orig. Error | Corrected Value | Error After Correction |
|--------|-------------|-------------|-----------------|------------------------|
| 1820   | 1823        | -3          | 1825.2          | 2.2                    |
| 2074   | 2073        | 1           | 2072.6          | -0.4                   |
| 2327   | 2322        | 5           | 2319.0          | -3.0                   |
| 2588   | 2573        | 15          | 2573.3          | 0.3                    |
| 2845   | 2823        | 22          | 2823.6          | 0.6                    |
| 3103   | 3077        | 26          | 3074.9          | -2.1                   |
| 3365   | 3329        | 36          | 3330.1          | 1.1                    |
| 3623   | 3581        | 42          | 3581.5          | 0.5                    |

Fitted correction function:  $f(d) = 0.9741d + 52.299$

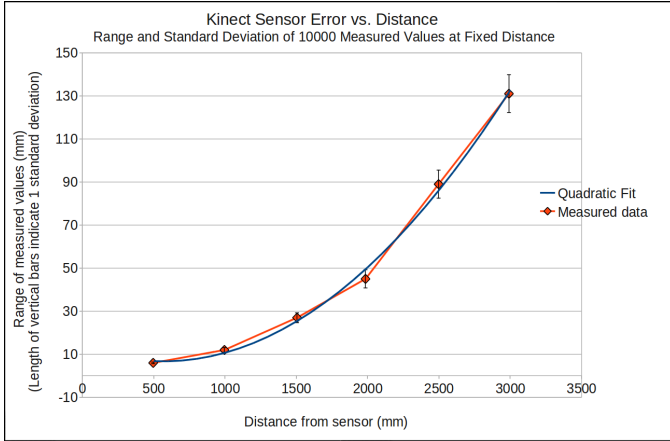


Figure 6: Kinect depth sensor precision with distance. Measured values show quadratic relationship between the distance to the depth camera and the range and standard deviation of depth values.

To verify this, we positioned a planar target parallel to the IR camera’s image plane and recorded a  $100 \times 100$  grid of depth measurements at the center of the depth image. We performed this experiment at distances of 0.5 m (device minimum range) to 3.0 m (beyond the maximum range used in our system) at intervals of 0.5 m.

Figure 6 shows the min-max range and standard deviation at each test distance from 0.5 m to 3.0 m, fitting closely to a quadratic falloff.

#### 4.2. Multi-Kinect Interference Problem and Solution

**The Multi-Kinect Interference Problem.** Since each Kinect unit projects the same dot pattern at the same wavelength, each Kinect unit is able to see the projected patterns of all other units and may have trouble distinguishing other units’ patterns from its own.

This problem is illustrated in Figure 7. A box was placed near the minimum depth range (0.6 m) of two Kinect units; their projected patterns overlap prominently and cause interference. Although significant interference is shown in the third column of the figure (there many small areas of missing data, or “holes”), we find two promising aspects of these results. First, the difference between the depth image with and without interference corresponds mostly to the missing data, not large differences in depth values; one needs primarily to fill missing points rather than correct grossly erroneous depth values (although some additional noise is also present). Second, one can see in the third column of the figure that the missing data varies between depth cameras – to some extent, redundant coverage between units allows depth cameras to fill in each other’s holes.

**Hardware Solutions.** We considered, but rejected several hardware solutions to the multi-Kinect interference problem. We contemplated installing a set of alternating synchronized shutters over each unit’s IR projector and camera so that each unit would see only its own dot pattern, as was explored in [18]. A serious disadvantage of this approach is that it would reduce

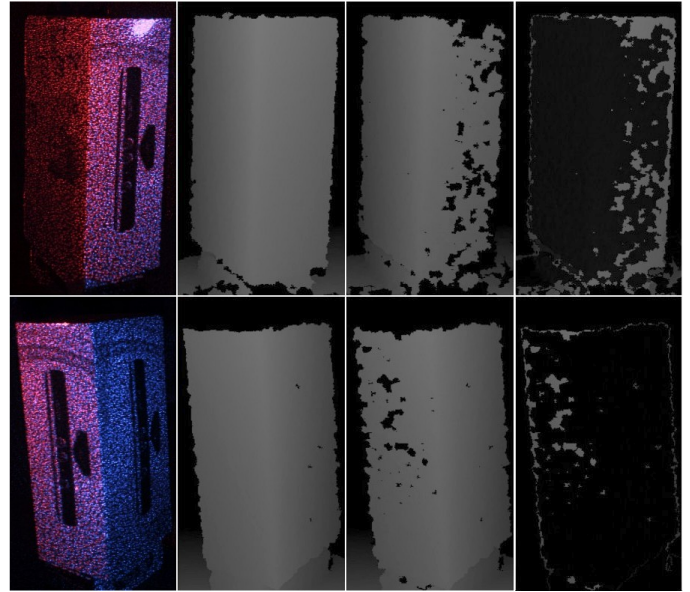


Figure 7: Kinect interference problem. First column: IR images showing combined projected dot pattern from camera 1 (red dots) and camera 2 (blue dots). Second column: depth images with no interference. Third column: depth images with interference from other camera. Fourth column: difference of second and third columns. Rows: data from each of two cameras.

frame rate or reduce the light available to the IR camera, depending on how the shutters are used. Another technique considered, but also ruled out, was IR filtering. We measured a group of eight Kinect units with a spectrometer and found that the peak-to-peak range of wavelengths was 2.6 nm, which we found too close to filter practically.

**Software Solutions.** As we did not find a suitable hardware solution to the Kinect interference problem, we looked to software solutions. As mentioned, it is fortunate that the Kinect generally returns no data rather than erroneous data when interference occurs. However, there are other situations in which the sensor returns no depth data. Due to the offset between the IR projector and camera, there are typically some surfaces “in shadow” that can be seen by the camera but receive no projected pattern from the IR laser due to occlusion. Additionally, surfaces may not be seen by the depth camera if they reflect little infrared light or are highly specular. An effective software solution should be able to fill small holes (making the assumption that they are part of a continuous surface), while ignoring large missing surfaces. We hope that the missing large surfaces are captured by another camera that observes the surface from a different location. Also, Kinect interference causes a small amount of high frequency depth noise that should be smoothed.

**Hole Filling.** We aimed for a solution that fills small holes and provides an initial depth smoothing, but leaves alone large missing surfaces – we do not want to make large assumptions about missing data in the scene. The obvious starting point for such an approach is a simple smoothing filter (such as Gaussian, box, median, or bilateral), but our application induces additional requirements:

- 371 1. Edge preservation: we do not want to introduce new depth  
 372 values across depth discontinuities – naive smoothing could  
 373 result in depth values floating in space. Additionally,  
 374 depth images are aligned to color textures, so color and  
 375 depth edges should coincide. A small depth edge shift  
 376 could cause a texture to be assigned to a physically dis-  
 377 tant surface. Depth edges at missing data boundaries  
 378 must be preserved or geometry may expand or contract.
- 379 2. Scale independence: from observation, depth noise ap-  
 380 pears at a higher spatial frequency than holes. Smoothing  
 381 should take place on a smaller scale than hole filling.

382 A standard median filter meets the first requirement (edge  
 383 preservation – although not pixel-exact). We devised a fast  
 384 modified median filter (Algorithm 1) that is effective at hole  
 385 filling while supporting the requirements above.

386 To allow for scale independence, a two-pass approach is  
 387 used. In the first pass, an expanded filtering window is used  
 388 to fill larger holes, but no smoothing is applied (i.e. only miss-  
 389 ing values are modified). In the second pass, a smaller win-  
 390 dow is used to fill any remaining small holes and provide initial  
 391 smoothing of the depth image. This method is similar to that  
 392 used in [19], but we use different criteria to determine when the  
 393 filter is applied.

394 To ensure that edges are preserved precisely and non-holes  
 395 are ignored, we apply three constraints to the filtering window:  
 396 a minimum amount of valid data must be present ( $t_c$ ), a mini-  
 397 mum amount of data must be present at the window edges ( $t_e$ ),  
 398 and the range of values in the window must be within a thresh-  
 399 old ( $t_r$ ). At each pixel, if the window constraints are not met, the  
 400 pixel is left unmodified. These thresholds and heuristics were  
 401 determined by applying a conventional median filter to sample  
 402 depth data and inspecting cases that did not meet our require-  
 403 ments listed above.

404 To enhance our filtering implementation in the previous sys-  
 405 tem [12], we incorporated the trimming operation of [17] into  
 406 our hole filling and initial smoothing filter to reduce the appear-  
 407 ance of ragged edges at depth discontinuities. This operation  
 408 rejects geometry that does not meet the previously described  
 409 data population and range tests, which was found empirically to  
 410 occur at object edges. This operation is performed only during  
 411 the second (small window) pass as to prevent over-trimming.

412 *GPU Implementation.* Our enhanced median filter implemen-  
 413 tation is based on a conventional median filter implementation  
 414 by McGuire [20], which uses a branchless hardcoded selection  
 415 algorithm to obtain the median for fixed radii. To provide high  
 416 performance for larger radii, we find the approximate median  
 417 by sampling over the filtering window. The median filter is  
 418 written as a fragment shader in the OpenGL Shading Language,  
 419 using textures to exchange data between passes.

### 420 4.3. Smoothing and Temporal Filtering

421 In our previous system [12], the smoothing described in  
 422 Section 4.2 was the sole filtering operation applied to the depth  
 423 maps. Although this filter made some improvement to the nois-  
 424 iness of the depth maps, several disadvantages were observed:

---

#### Algorithm 1 Modified Two-Pass Median Filter for Hole Filling

---

```

for  $pass = 1$  to  $2$  do
  for  $i = 1$  to  $numPixels$  do
     $depth\_out[i] \leftarrow depth\_in[i]$ 
    if  $depth\_in[i] = 0$  or  $pass = 2$  then
      {Perform filtering tests}
       $count \leftarrow 0, enclosed \leftarrow 0$ 
       $v \leftarrow \{\}, n \leftarrow neighbors(depth\_in[i], radius_{pass})$ 
       $min \leftarrow min(n), max \leftarrow max(n)$ 
      for  $j = 1$  to  $n.length$  do
        if  $n[j] \neq 0$  then
           $count \leftarrow count + 1$ 
           $v[count] \leftarrow n[j]$ 
          if  $on\_edge(j)$  then
             $enclosed \leftarrow enclosed + 1$ 
          end if
        end if
      end for
      if  $max - min \leq t_r$ , and  $count \geq t_c$  and  $enclosed \geq t_e$  then
        {If filtering tests passed, find median}
         $sort(v)$ 
         $depth\_out[i] \leftarrow v[v.length/2]$ 
      else if  $pass = 2$  then
        {If filtering tests failed on 2nd pass, trim}
         $depth\_out[i] \leftarrow 0$ 
      end if
    end if
  end for
if  $pass = 1$  then
   $depth\_in \leftarrow depth\_out$ 
end if
end for

```

---

- 425 1. The median filter does not introduce new depth values  
 426 while smoothing, and at greater distances, the steps in  
 427 depth returned by the Kinect become large. This can re-  
 428 sult in significant depth discontinuities to appear among  
 429 smoothed regions that are physically continuous.
- 430 2. No temporal noise suppression is performed, causing sur-  
 431 faces to have a jittering effect.

To address item 1 above, we further smoothed our hole  
 filled depth maps with the following filter, based on the bilateral  
 filter of [21]:

$$\frac{1}{W_{x,y}} \sum_{u=-r}^r \sum_{v=-r}^r [D(x+u, y+v) g(\|u, v\|_2, \sigma_r) g(|D(x+u, y+v) - D(x, y)|, \sigma_d)]$$

432 where  $(x, y)$  is the coordinates of the pixel to be filtered,  $r$  is  
 433 the filter radius,  $D(u, v)$  is the depth at location  $(u, v)$ ,  $g(t, \sigma)$   
 434 is the Gaussian function  $e^{-t^2/\sigma^2}$ ,  $\sigma_r$  and  $\sigma_d$  control the Gaussian  
 435 falloff across the filtering window and depth respectively, and  
 436  $W_{x,y}$  is a normalizing constant equal to the sum of the weights.  
 437 This filter introduces new depth values while smoothing, but



Figure 8: Fast Meshing for the Kinect. Left: Triangle mesh template stored in GPU memory. Center-Left: Vertex shader extrudes template using camera intrinsics and depth texture. Center-Right: Geometry shader rejects triangles corresponding to missing or discontinuous surfaces (shown in red). Right: Resultant textured mesh.

438 weighted values fall off sharply with increased depth distance  
 439 from the central element to prevent smoothing over depth dis-  
 440 continuities.

441 As a first attempt to address item 2 above, we performed  
 442 simple temporal noise suppression by updating values in our  
 443 depth map only if they exceed a threshold from the last recorded  
 444 value. The depth-dependent threshold was set conservatively  
 445 based on the measured noise values plotted in Figure 6.

446 *GPU Implementation.* The bilateral filter was implemented as  
 447 an OpenGL fragment shader which takes the hole-filled and  
 448 pre-smoothed depth maps of Section 4.2 as input and produces  
 449 a final depth map texture to use for mesh generation. The tem-  
 450 poral filter was incorporated as an initial threshold check to the  
 451 bilateral filter. The final depth map is copied within GPU mem-  
 452 ory for comparison with the next frame.

#### 453 4.4. Mesh Generation

454 The Kinect provides per-pixel depth readings that are gen-  
 455 erally too sparse to render directly as small fixed-size points.  
 456 Therefore it is useful to create a surface representation using  
 457 the depth data. Our requirements for surface generation are as  
 458 follows:

- 459 1. Must be continuous if physical surface is continuous
- 460 2. Must work in situations with missing data, as is common  
 461 with Kinect
- 462 3. Must detect and preserve depth discontinuities at edges
- 463 4. Must be fast (5 Kinects generate >45M depth readings/sec)

464 Although approaches exist [22] for directly rendering points  
 465 from multiple depth images, we chose a triangle mesh surface  
 466 representation as it meets these requirements and is also sup-  
 467 ported natively by graphics hardware. We use a simple mesh-  
 468 ing technique which is described in Algorithm 2 and illustrated  
 469 in Figure 8. The depth values from the Kinect sensor are used  
 470 to extrude vertices from a template triangle mesh. Any trian-  
 471 gle associated with a vertex that corresponds to a missing depth  
 472 value is rejected, as are those with a pair of vertices that ex-  
 473 ceeds a maximum depth threshold. Since the Kinect provides  
 474 depth data that varies in accuracy with depth, our depth thresh-  
 475 old varies with depth as well.

---

#### Algorithm 2 Mesh generation algorithm

---

```

for each candidate triangle do
   $t \leftarrow thresh_{depth\_discontinuity} + f_{depth\_err}(\min(depth_{v_i})) +$ 
   $f_{depth\_err}(\max(depth_{v_i}))$ 

  {Transform vertices from normalized image coordinates to
  camera coordinates in physical units}
  if  $depth_{v_i} \neq 0$  and  $abs(depth_{v_i} - depth_{v_j}) \leq t$ , for  $j > i$ 
  then
     $v_{i_x} \leftarrow \frac{(v_{i_x} - center\_proj_x)depth_{v_i}}{focal_x}$ 
     $v_{i_y} \leftarrow \frac{(v_{i_y} - center\_proj_y)depth_{v_i}}{focal_y}$ 
     $v_{i_z} \leftarrow depth_{v_i}$ 
  else
    reject triangle
  end if
end for

```

---

476 *GPU Implementation.* Our implementation of the simple mesh  
 477 algorithm takes advantage of the connectivity of a depth image,  
 478 requiring no geometry to be transferred to the GPU after pro-  
 479 gram initialization. At program start we generate a triangulated  
 480 plane at the Kinect’s depth resolution and store it in GPU mem-  
 481 ory. For each new frame, a vertex shader shapes the template  
 482 plane using camera intrinsics and Kinect depth values, which  
 483 are accessed through a texture map. A geometry shader is used  
 484 to reject triangles corresponding to missing depth values or dis-  
 485 continuous surfaces as described in Algorithm 2.

486 This approach is very bandwidth-efficient – it requires only  
 487 16 bits of depth information for each pair of triangles generated  
 488 and uses the depth map already transferred to GPU memory for  
 489 the hole filling process. The approach is also fast as all vertex  
 490 positions are generated on the GPU in parallel.

#### 491 4.5. Data Merger

492 *Overview.* A goal of our system is to provide coverage of all  
 493 surfaces that can be seen from the perspective of a remote user.  
 494 A single Kinect is not able to provide adequate coverage and  
 495 therefore a means to merge data between multiple units is nec-  
 496 essary. When generating meshes we did not discuss a means  
 497 to merge overlapping surfaces geometrically. Approaches used  
 498 in stereo vision, such as the visibility-based depth image fu-  
 499 sion method of Merrell et al. [19], generally assume high levels  
 500 of error and inconsistencies (outliers) between maps that must  
 501 be resolved. The Kinect’s structured-light based depth read-  
 502 ings, however, are generally free of such outliers and have a  
 503 low error at near range. In our application, Kinect sensors are  
 504 used at close proximity and we expect lower and predictable  
 505 error based on the angle and distance to the camera and mea-  
 506 sured calibration error. Therefore, we assume that the surfaces  
 507 have enough fidelity that we can simply draw them on top of  
 508 each other, avoiding the need for a geometric merger algorithm.  
 509 This has several performance advantages: the computational  
 510 expense of performing the merge is spared, runtime varies lin-  
 511 early with the number of cameras, surfaces for all cameras can



512 be processed in parallel, and a fast mesh generation technique  
513 (Section 4.4) can be used.

514 However, even though geometry is sufficiently accurate for  
515 our purposes, texture image quality may be poor. Z-fighting be-  
516 tween overlapping surfaces with textures that vary in resolution,  
517 color imbalances, and misalignment yields unpleasing results.  
518 Ideally, we want to utilize only the data from the camera with  
519 the highest resolution depth and color information available at  
520 a given surface, with a seamless transition to data from adjacent  
521 cameras.

522 Our approach addresses the problem of data merger in im-  
523 age space using a visibility-based approach. The data from each  
524 camera is rendered independently for the desired viewpoint, and  
525 color information is saved along with a depth and a quality es-  
526 timate at each pixel. When renderings for all cameras are com-  
527 plete, the depth values are used to determine which cameras can  
528 see the front surface. At each pixel, the color values of cameras  
529 with a view of the front surface are weighted by the quality es-  
530 timates. This process is illustrated in Figure 9.

531 *Texture Quality and Depth Error Estimation.* Since our approach  
532 relies on the notion of a “quality” measurement at each pixel,  
533 we provide an estimate based on resolution – the area on the im-  
534 age sensor available to determine the pixel’s color or position.  
535 The area is estimated using the cosine of the angle between the  
536 surface normal of the pixel and the squared distance from the  
537 pixel to the image sensor. The relationship between area and  
538 resolution is straightforward for a color image, and we demon-  
539 strated previously that the Kinect depth error increases quadrat-  
540 ically. We approximate quality by assuming that both color and  
541 depth error increase quadratically, yielding the quality value in  
542 Equation 1.

$$quality = \left( \frac{\cos\theta_{normal \rightarrow camera}}{distance^2} \right)^2 \quad (1)$$

543 Note that this formulation is similar to a diffuse lighting cal-  
544 culation with attenuation (for a light positioned at the sensor’s  
545 location) that can rapidly be performed on almost any graphics  
546 hardware.

Our approach also requires determination of which pixels  
represent the closest surface with respect to viewing position.  
We store the depth values at each pixel, but due to calibration  
and depth sensor error the values corresponding to the front sur-  
face do not coincide exactly. Equation 2 is used to estimate the  
range of each depth position, so that the range of depths corre-  
sponding to the front surface can be determined.

$$[-err_{calib} - f_{depth\_err}(depth), err_{calib} + f_{depth\_err}(depth)] \quad (2)$$

547 Calibration error ( $err_{calib}$ ) can be estimated using the re-  
548 projection error that is returned by the camera calibration rou-  
549 tine. Depth error ( $f_{depth\_err}$ ) can be estimated using the data  
550 from Figure 6.

551 *Data Merger Algorithm.* Algorithm 3 describes the process of  
552 merging the renderings for each camera. At each pixel, the front  
553 surface tolerance is determined by finding the closest depth

554 value that represents the far end of any pixel’s estimated depth  
555 range. The color values for all pixels with this depth value or  
556 nearer are weighted by quality to obtain the final pixel color.

---

**Algorithm 3** Data merger algorithm

---

```

for each output pixel p do
  depthfar ← ∞
  {Determine far bound of closest surface}
  for each camera c do
    dfar ← depthcp + errcalib + fdepth_err(depthcp)
    if dfar < depthfar then
      depthfar ← dfar
    end if
  end for
  colorsum ← 0, qualitysum ← 0
  for each camera c do
    {Only consider cameras with view of closest surface}
    if depthcp ≤ depthfar then
      matchbest ← ∞
      {Perform photometric search for closest matching
      pixel to camera with best quality estimate}
      for each pixel s in search window do
        match ← ||colorqp - colorcs||2
        if match < matchbest then
          matchbest ← match
          colorbest ← colorcs
          qualitybest ← qualitycs
        end if
      end for
      {If photometric threshold met, perform quality
      weighted blending}
      if matchbest ≤ threshmatch then
        colorsum ← colorsum + qualitybest * colorbest
        qualitysum ← qualitysum + qualitybest
      end if
    end if
  end for
  coloroutput ← colorsum / qualitysum
end for

```

---

557 *Photometric Constraint.* A shortcoming of our original data  
558 merger algorithm [12] was the unconditional blending of data  
559 between all cameras regardless of color consistency. If intensity  
560 varied widely between cameras (due to geometric misalignment  
561 or specularities), the resultant combined textures would appear  
562 blurry or a double image would appear at conflicting pixels.  
563 Inspired by [7], we apply a photometric search and constraint  
564 before blending each quality weighted color sample. For each  
565 camera, we perform a search within a small window for the  
566 closest matching color to the color associated with the camera  
567 with the highest quality estimate ( $color_{q_p}$ ) and blend only if the  
568 colors match within a threshold.

569 *GPU Implementation.* Our fast GPU implementation supports  
570 calculation of depth, quality, and color values in one pass per  
571 camera and allows all cameras’ renderings to be merged at once

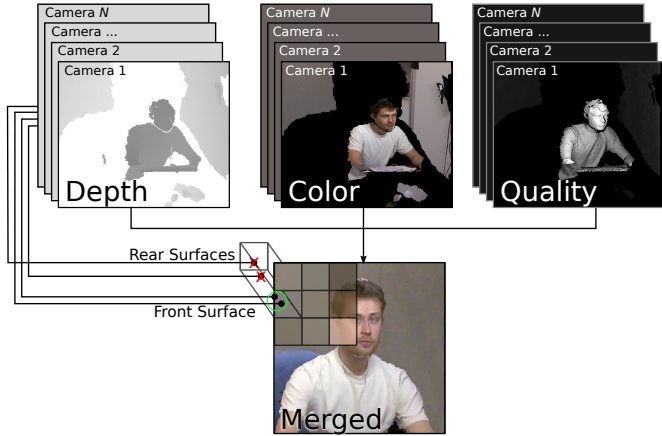


Figure 9: Data Merger Algorithm. Depth, color, and quality estimate values are determined at each pixel for each camera. The front surface is determined using the depth information, and the associated color values are weighted by the quality estimates.

572 in a second pass. When generating the triangle mesh in an  
 573 OpenGL geometry shader, we compute the distance to the camera  
 574 and the angle between the camera and surface normal and  
 575 save these values as vertex attributes. During rasterization, an  
 576 OpenGL fragment shader computes a color value and a quality  
 577 value (using Equation 1) at each pixel, storing the quality value  
 578 in the alpha channel and the depth value from the Z-buffer in a  
 579 separate texture. When the renderings for all cameras are complete,  
 580 all data is merged in an OpenGL fragment shader according to  
 581 Algorithm 3.

#### 582 4.6. Multiple Camera Color Matching

583 *Overview.* The need for color matching is common for many  
 584 camera systems, as even the same model device may exhibit  
 585 different color gamuts [23]. This need is exacerbated in inex-  
 586 pensive devices like the Kinect sensor, which allows only auto-  
 587 matic color and exposure control (with present drivers), yield-  
 588 ing color values that may vary dramatically between adjacent  
 589 cameras. Here traditional color matching techniques, such as  
 590 adjusting color to match a physical target seen by each cam-  
 591 era, are ineffective because automatic control may alter color  
 592 balances at any time. We present an automatic color matching  
 593 technique that uses depth information to find color correspon-  
 594 dences between cameras, which can be used to build a color  
 595 matching function. We believe this technique may be useful  
 596 when manual color adjustment is unavailable, or as a fast ap-  
 597 proximate alternative to conventional matching techniques.

598 *Obtaining Color Correspondences.* To build a set of color cor-  
 599 respondences between cameras, we first find pairs of points  
 600 from two cameras that correspond to approximately the same  
 601 point in 3D space. We assume that each pair of points repre-  
 602 sents the same point on a diffuse surface in physical space, and  
 603 therefore should agree in color. To find these point correspon-  
 604 dences, we refer to our previously described visibility-based  
 605 data merger algorithm. The algorithm rendered the scene in-  
 606 dividually for each Kinect camera and examined corresponding

607 depth values to determine which represented the front surface.  
 608 For color matching, if two cameras have depth values that repre-  
 609 sent the front surface at a given pixel, we add their color values  
 610 to a list of correspondences.

611 Since this approach is visibility-based, the color correspon-  
 612 dences obtained are sensitive to the position of the virtual cam-  
 613 era. If the same virtual camera position is used for color match-  
 614 ing and rendering, color matching is tailored to the colors actu-  
 615 ally seen by the user. However, if a pair of cameras have few  
 616 surfaces in common from the viewpoint used for rendering, or  
 617 if these surfaces have a limited range of colors, there may be  
 618 too few correspondences to build a robust color matching func-  
 619 tion. In this case, point correspondences can be computed from  
 620 a reference view (such as a bird’s eye view), rather than from  
 621 the view used for rendering. To build more robust color corre-  
 622 spondences, additional techniques could be used. For example,  
 623 the color correspondences could be built from renderings from  
 624 several viewpoints, or could be collected over time.

625 *Building a Color Matching Function.* There are many advanced  
 626 techniques for building color matching functions from a set of  
 627 color correspondences, such as that of Ilie and Welsh [23]. To  
 628 demonstrate our approach, we used a simple method – color  
 629 correspondences were fit to a linear model. Since our color cor-  
 630 respondences were noisy (small errors in surface position may  
 631 result in a large difference in color), we used the RANSAC [24]  
 632 method for fitting, which is robust to outliers. Figure 10 shows  
 633 a plot of actual color correspondences (for one channel) and the  
 634 fitted linear color matching function.

635 *Implementation.* For our test setup, we matched the colors of  
 636 each camera to our ceiling-mounted master camera. We elected  
 637 not to run the color matching function on every frame, as small  
 638 variations in color matching functions resulted in a color cyc-  
 639 cling effect. Instead a new color matching function was built  
 640 whenever the user pressed a function key. As real-time perfor-  
 641 mance was not needed, we implemented color matching func-  
 642 tionality on the CPU. We believe our implementation could  
 643 be improved by performing bundle adjustment across cameras  
 644 and by running the color matching function automatically when  
 645 some criteria is met.

#### 646 4.7. Eye Position Tracking

647 *Overview.* To allow our system to render a set of correct stereo  
 648 viewpoints from the user’s position, we need to obtain the po-  
 649 sition of the viewer’s eyes in 3D space. Many approaches to  
 650 tracking have been devised, such as measuring the magnetic  
 651 field around a marker, segmenting and triangulating the posi-  
 652 tion of reflective markers as seen by an array of cameras, and  
 653 using computer vision techniques to recognize objects in im-  
 654 ages. The latter approach has been used to obtain the 3D posi-  
 655 tions of eyes with a conventional 2D camera, but assumptions  
 656 or measurements must be made of the face. We aim to improve  
 657 these techniques by incorporating depth information. One im-  
 658 pressive recent approach [25] used depth information to build a  
 659 deformable mesh that was tracked to a user’s face in real-time,  
 660 but required a 6.7 second initialization time and achieved only

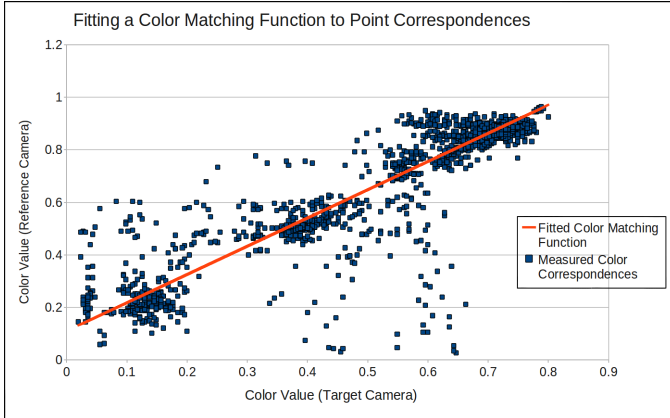


Figure 10: Color matching using 3D point correspondences. Plot shows color correspondences between a pair of cameras for one color channel and the RANSAC-fitted linear color matching function.

661 moderate real-time performance (10-12 Hz). Since we require  
 662 higher performance and do not need tracking of the entire face,  
 663 we look to an alternate approach – performing 2D eye detec-  
 664 tion and transforming the detected position into 3D using the  
 665 Kinect’s depth data.

666 *Requirements.* Our tracking system should meet the following  
 667 requirements for use in our telepresence system:

- 668 1. Accuracy: at a 1 m distance, 15 mm of lateral movement  
 669 causes the eye to sweep over one display subpixel seen  
 670 through the barrier of our autostereo display; for best  
 671 quality results tracking accuracy should be  $\pm 7.5$ mm.
- 672 2. Speed, Latency: we do not anticipate rapid head move-  
 673 ments in our application. To support the modest move-  
 674 ment of 25 cm/sec, framerate must be  $> 33.3$  Hz and la-  
 675 tency must be  $< 30$ ms to meet the accuracy requirements  
 676 above.

677 *2D Eye Tracking.* To perform 2D eye detection on the color im-  
 678 age, we use Viola [26] and Lienhart’s [27] approach of boosted  
 679 Haar classifiers, as implemented in OpenCV. First the face is  
 680 detected (using a classifier from Lienhart), and then eyes are de-  
 681 tected in the facial region (using classifiers from Castrillon [28]).  
 682 Once the eyes are found, their pattern is saved and subsequent  
 683 eye searches are performed using normalized cross correlation.  
 684 An image pyramid is used to accelerate the cross correlation  
 685 search. If the strongest response to cross correlation falls be-  
 686 low a threshold, detectors are again used to locate facial fea-  
 687 tures. The face is first searched for in the region surrounding  
 688 the last known eye position; if not found the entire image is  
 689 again searched. All detection and tracking operations were per-  
 690 formed in the CPU, as it was not heavily utilized elsewhere in  
 691 our system. A single Kinect unit, mounted above the autostereo  
 692 display, was used for tracking.

693 *Using Depth to Obtain 3D Eye Position.* Once the center of  
 694 both eyes have been detected, the 2D position is transformed  
 695 into 3D using the Kinect’s depth information and measured

696 camera intrinsics and extrinsics. To reduce the effects of noise  
 697 and missing data, depth values are averaged over a small radius  
 698 around the eye position. A Kalman filter was used to improve  
 699 the accuracy and stability of the 3D tracked eye positions as  
 700 well as predict the locations of the eyes between sensor read-  
 701 ings. Although our tracking system requires no prior measure-  
 702 ments of the user’s face, accuracy can be improved if the true  
 703 interpupillary distance (IPD) is known. If the system is utilized  
 704 by a single user over a capture session, an accurate IPD estimate  
 705 can be learned over time.

706 *Discussion.* Our tracking system offers several advantages over  
 707 existing systems. It uses inexpensive hardware (the Kinect sen-  
 708 sor) and allows the same device to be used for both tracking and  
 709 3D capture at the same time. Since the eyes are tracked  
 710 independently, our system allows correct calculation of 3D eye  
 711 positions without measurements or assumptions of face size or  
 712 IPD.

713 We believe our system could be improved with a more ro-  
 714 bust set of feature detectors – our current system allows for only  
 715 moderate head rotations and does not work well with glasses.  
 716 Depth data could also be further utilized to improve speed; for  
 717 example, the face search area could be restricted to depths that  
 718 are within the range of a seated user. Multiple cameras could be  
 719 utilized to offer better coverage of a rotated head or to improve  
 720 the accuracy of the system.

#### 721 4.8. Stereo Display

722 *Overview.* As mentioned, research [13] has shown that stereo  
 723 displays can increase the sense of shared presence, although  
 724 systems requiring 3D glasses obstruct eye contact and have  
 725 been found to be disruptive to most users. Therefore, we desire  
 726 an autostereo display for our system.

727 *Display Selection.* Our display system should meet the follow-  
 728 ing requirements for use in our telepresence system:

- 729 1. Preservation of full captured color and detail.
- 730 2. Large enough to allow remote scene to be observed as  
 731 life-sized at proper viewing distance.
- 732 3. Support for continuous viewpoints and horizontal and ver-  
 733 tical parallax.
- 734 4. Support for a range of movement typical of a seated user.
- 735 5. Interactive update rates that meet our tracking require-  
 736 ments.

737 We were in possession of a fixed parallax barrier display  
 738 that met these requirements – an X3D-40 display by X3D tech-  
 739 nologies (circa 2004). The display measures 40 in diagonally  
 740 and has a 1280×768 pixel resolution and a 60 Hz update rate.  
 741 Since the display supports only a limited number of views, track-  
 742 ing was employed. In a future system, we intend to utilize a  
 743 display that supports multiple users, such as the Random Hole  
 744 display of Ye et al [29].

745 *Rendering for the Display.* Since our system uses head track-  
746 ing, we rendered views for the display using off-axis frustra  
747 between the eyes and the display. The position of the eyes  
748 was determined using the tracking system, and the position of  
749 the monitor was measured using our 3D capture system. An  
750 OpenGL fragment shader was used to generate the diagonally  
751 interleaved pattern needed by our parallax barrier display for  
752 each pair of stereo views.

753 *Tracking Problem and Intra-Frame Rendering.* While using our  
754 fixed parallax barrier display, a user may see an incorrect view  
755 or significant artifacts (dark black bands or fuzziness) if out of  
756 the expected viewing position. If the rendering update rate is  
757 lower than the rate required by our tracking system, a user may  
758 experience these effects if moving, resulting in poor stereo per-  
759 ception.

760 This problem has been addressed previously for dynamic  
761 barrier displays [30] by generating the parallax barrier *stripes*  
762 asynchronously at higher rates than rendering takes place. For  
763 fixed barrier displays, we developed a new technique to address  
764 this problem – rendering barrier *display patterns* at a higher rate  
765 while new frames are rendered more slowly offscreen.

766 Since the time it takes to generate a parallax barrier pattern  
767 for a new eye position is very short and fixed with our GPU  
768 implementation, we can draw one or more new barrier patterns  
769 while in the process of rendering a frame for the next view-  
770 ing perspective. These intra-frame barrier patterns use the new  
771 estimated eye position and the last rendered viewing position,  
772 saved in textures. Using OpenGL, we are able to draw to the  
773 screen mid-frame by switching between multiple frame buffers.  
774 To keep our parallax barrier generation rate and rendering rate  
775 independent, we stop to draw a new barrier pattern whenever  
776 a fixed amount of time has elapsed during rendering, periodi-  
777 cally flushing the pipeline to allow for better time granularity  
778 between asynchronous GL calls. The result is a high fixed bar-  
779 rier display rate, independent of rendering rate, at the expense  
780 of a small decrease in rendering rate. Specific rates are listed in  
781 Section 5.4

## 782 5. Results

### 783 5.1. Camera Coverage and Calibration Results

784 *Camera Coverage.* Our camera arrangement (shown in upper  
785 left of Figure 3), includes most of the surfaces seen by a seated  
786 remote user, as shown in Figure 11. Overlapping camera cov-  
787 erage preserves large surfaces on the rear wall, which would  
788 otherwise be occluded by the seated user. Redundant coverage  
789 also helps prevent self shadowing. For example, in Figure 1 the  
790 coffee cup occludes part of the user’s chest with respect to one  
791 of the front cameras, but the missing surfaces are filled with  
792 data from the other front camera.

793 *Depth Bias Correction.* To test the depth bias correction func-  
794 tions determined in Section 4.1, two checkerboard targets were  
795 placed centered in front of two Kinect units, so that nearer and  
796 farther targets were approximately 3.0 m and 3.5 m from the  
797 units respectively. The Kinects were aimed at the targets and

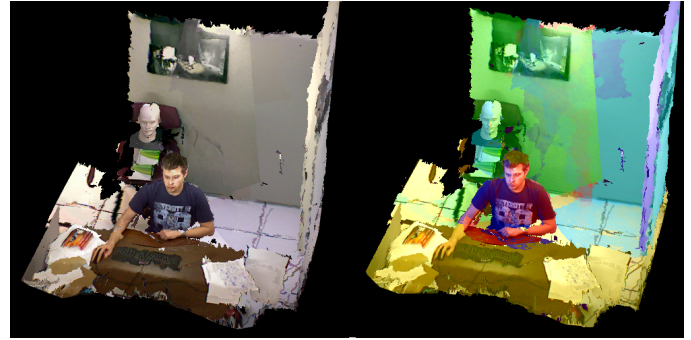


Figure 11: Camera Coverage. Left: Camera coverage in our cubicle area using five depth cameras. Right: Color coded contributions of individual cameras.

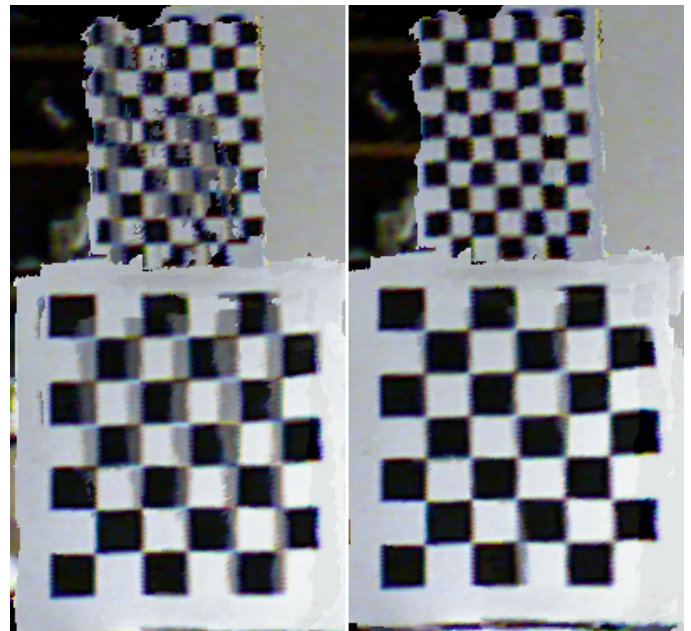


Figure 12: Depth bias correction results. Left: No depth bias correction applied. Right: Depth bias correction applied.

798 positioned with a wide baseline of 2.5 m to accentuate errors  
799 in depth. The data from the two cameras was rendered using  
800 the techniques described in Section 4, except that the photo-  
801 metric constraint of Section 4.5 was disabled during blending  
802 so that the misalignments of the high contrast targets could be  
803 observed.

804 The results are shown in Figure 12. One can see a signif-  
805 icant improvement in the alignment of both checkerboards be-  
806 tween cameras, although a small amount of misalignment re-  
807 mains.

### 808 5.2. Data Processing and Rendering Results

809 *Mesh Generation.* All images in Figure 13 show the result of  
810 mesh generation. Our requirements are met: the mesh offers  
811 a continuous surface, discontinuous surfaces (such as from the  
812 body to the rear wall) are properly separated, and missing data  
813 (small area under chin) is tolerated.

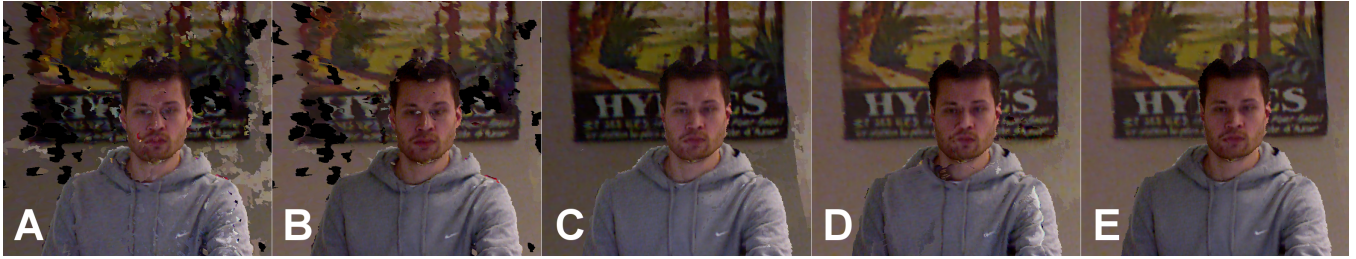


Figure 13: Data processing results. A: No enhancements applied. B: All enhancements except hole filling and smoothing. C: All enhancements except color matching. D: All enhancements except quality weighted data merger. E: All enhancements applied. (“All enhancements” includes to hole filling, smoothing, data merger, and color matching)



Figure 15: Data merger results. A1: No merger applied, meshes drawn on top of each other. A2: Merger with simple average. A3: Merger with quality weighting. B1: Color coded camera contributions with no merger. B2: Camera contributions with quality-weighted merger.

814 *Hole Filling and Smoothing.* Image E of Figure 13 (as compared to image B of Figure 13) shows the result of the hole filling and smoothing filters applied on a scene with four overlapping cameras. In this example, 100% of holes caused by interference were filled while textures remained aligned to the mesh (rates of > 90% are typical).

820 The effect of the additional bilateral smoothing filter is illustrated in Figure 14. Note that the more distant surfaces (ex: wall) appear flatter and smoother with the addition of the bilateral filter (right image) than with the median-only smoothing of our previous system [12] (center image).

825 *Data Merger.* Image E of Figure 13 (as compared to image D of Figure 13) and all of Figure 15 show the result of the data merger algorithm on four cameras, which is cleaner and smoother than meshes simply drawn over each other or averaged. In image B1 of Figure 15, one can see in the unmerged example that the mesh of the right-front camera (tinted blue) is drawn entirely over the data from the left-front camera (tinted red). These surfaces should coincide exactly, but a small calibration error places the surface from right-front camera closer to the viewer. In image B2 of Figure 15, one can see that the quality-weighted merger algorithm smoothly transitions between camera data across the face.

837 Figure 16 shows the improvement gained when the photometric constraint is incorporated into the data merger algorithm. Without the photometric constraint, some areas of the poster (for example, around the walking figures and the large letter È) have ghost images due to slightly misaligned surfaces between cameras. With the constraint enabled, the ghost images have mostly disappeared, although a faint outline is still visible.



Figure 16: Photometric constraint. Top: No photometric constraint in data merger algorithm. Bottom: Photometric constraint applied.



Figure 14: Smoothing results. Left: No smoothing (or hole filling) applied. Center: Median smoothing. Right: Bilateral smoothing.



Figure 17: Cumulative enhancements. Left: Rendering using our previous system [12]. Right: Rendering using the improved calibration, data processing, and rendering techniques described in this paper.

Table 2: Tracking performance over 1600 frames.

| Case                                       | #    | %     | Avg Time(ms) |
|--|------|-------|--------------|
| Eyes found (full frame face/eye detect)    | 3    | 0.19  | 140.5        |
| Eyes found (partial frame face/eye detect) | 9    | 0.56  | 19.8         |
| Eyes found (pattern search)                | 1585 | 99.06 | 2.3          |
| Eyes not found                             | 3    | 0.19  | 13.6         |

844 *Color Matching.* Image E of Figure 13 (as compared to im-  
 845 age C of Figure 13) shows the result of color matching in a  
 846 scene with four cameras. The unmodified image shows small  
 847 color inconsistencies between the front cameras (on the face  
 848 and clothing) and significant inconsistencies between the four  
 849 cameras that overlap in the background – most notably on the  
 850 far right side of the image. The automatic color matching al-  
 851 gorithm mostly resolved these deviations, although some faint  
 852 color seams are still visible.

853 *Cumulative Result.* Figure 17 shows the cumulative effect of  
 854 the enhanced data processing and rendering (as well as calibra-  
 855 tion) techniques over those from our previous system [12]. Note  
 856 the improved straightness of lines on the poster and checker-  
 857 board, clearer textures on the face, t-shirt, and poster, and the  
 858 absence of ghosting on the checkerboard.

### 859 5.3. Eye Tracking and Stereo Display Results

860 *Eye Detection Rate and Speed.* Table 2 shows the tracking per-  
 861 formance typical of a seated user over a 1600 frame sequence.  
 862 For the sequence, the user was seated centered 1 m from the  
 863 display and tracking camera and moved his head left, right, for-  
 864 ward and backward over a range of  $\pm 0.5$  m. The average head  
 865 movement speed was 48 cm/s, measured using the detected 3D  
 866 eye positions. Positive eye detection occurred on  $>99\%$  of the  
 867 frames at an average rate of 2.7 ms. In the worst case, when the  
 868 face was lost and the entire frame had to be searched, a notice-  
 869 able delay of 140.5 ms on average occurred.

870 *Tracking Performance.* Figure 18 provides a measure of the  
 871 performance of the eye tracking by comparing the 3D distance  
 872 between a pair of tracked eyes and the true measured interpup-  
 873 ilarly distance (IPD). IPD was used as the ground truth for accu-  
 874 racy as we were not in possession of equipment that would al-  
 875 low us to measure our positional accuracy directly. This metric  
 876 was measured over a sequence of 1761 frames, in which a user  
 877 seated 1 m from the tracking camera moved his head to the left,  
 878 right, forward and backward over a range of  $\pm 0.5$  m. 85.6% of  
 879 measurements were within  $\pm 5$  mm of the true IPD, and 96.4%  
 880 were within  $\pm 10$  mm.

881 *Tracking Accuracy and Stereo Quality.* Since our tracking sys-  
 882 tem is designed to support a stereo display, it is useful to test the  
 883 two systems together. To demonstrate that our tracking system  
 884 is fast and accurate enough to support our parallax barrier au-  
 885 tostereo display with good quality, we shot video of our system  
 886 through a tracking target (shown in Figure 20). Our tracking  
 887 system is able to detect the target as if it were a real face and

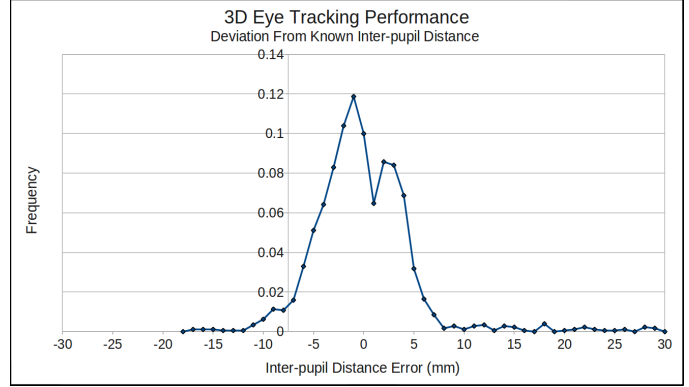


Figure 18: 3D Eye tracking performance. Plot shows measured deviations from a known inter-pupil distance.

888 thus one of the stereo views will be generated from the cor-  
 889 rect perspective of the camera placed behind an eye. Using this  
 890 setup, the target and camera were positioned 1.25 m from the  
 891 tracking camera and display and were moved at a rate of ap-  
 892 proximately 24 cm/sec.

893 Without tracking prediction and intra-frame rendering en-  
 894 abled, the rendering and parallax barrier pattern generation rate  
 895 was 21 Hz in our four camera test setup. As seen in the left of  
 896 Figure 19, results were very poor; the tracking and rendering  
 897 could not keep up with the target as it moved into the view-  
 898 ing zone intended for the other eye and thus both views could  
 899 be seen prominently and simultaneously. With tracking predic-  
 900 tion and intra-frame rendering enabled (right of Figure 19), the  
 901 rendering rate dropped slightly to 18 Hz but the barrier genera-  
 902 tion rate more than doubled to 48 Hz. Results were much im-  
 903 proved – the view seen by the camera is crisp and only very faint  
 904 ghosting can be seen to the right of the mannequin head and  
 905 box. [Please note that the performance numbers quoted in this  
 906 paragraph and the images of Figure 19 were measured using  
 907 our previous system [12] and performance has since increased  
 908 due to a graphics card upgrade (see Table 3). Although the  
 909 higher display rates of our upgraded hardware lessen the need  
 910 for intra-frame rendering in our current setup, we believe the  
 911 technique remains applicable to those with more modest hard-  
 912 ware. It will also allow us to maintain stereo quality as more  
 913 advanced (and computationally complex) reconstruction algo-  
 914 rithms are utilized and more Kinects are added to the system in  
 915 the future, and will allow increased head movement speeds.]

### 916 5.4. System Performance

917 Table 3 lists the performance achieved with our test system  
 918 in various configurations. When rendering for a single view, the  
 919 system was able to maintain average frame rates of 48 Hz for  
 920 five depth cameras with all enhancements (meshing, hole fill-  
 921 ing, quality-weighted data merger) enabled. For tracked stereo  
 922 configurations, rendering rates fell to 34 Hz, but a parallax bar-  
 923 rier pattern rate of 74 Hz preserves smooth head tracking and  
 924 stereo quality.

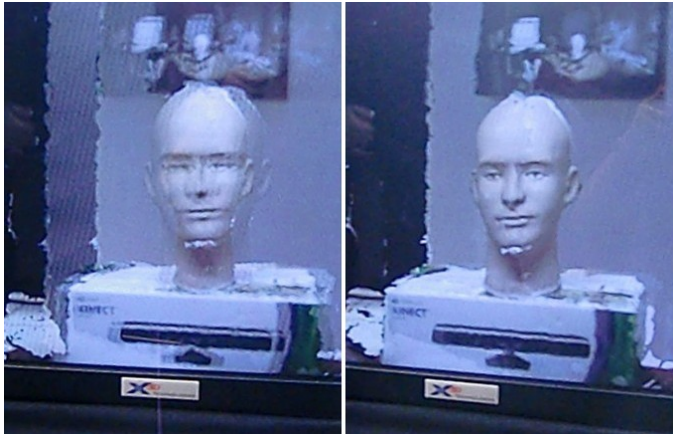


Figure 19: Head-tracked stereo in motion. Left: Tracking prediction, intra-frame rendering disabled. Right: Prediction, intra-frame rendering enabled. (Note: faint image on right side is reflection of photographer).

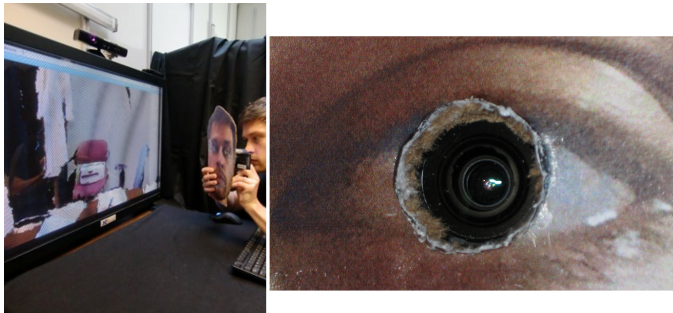


Figure 20: Tracking target. Left: Head cutout used to test eye tracking, with camera behind eye. Right: close-up of camera through eye.

Table 3: Average display rates (frames per second) for 5 Kinects

| Single View w/ Selected Enhancements         |     |
|--|-----|
| Meshing                                      | 160 |
| Meshing, Hole Filling/Smoothing              | 117 |
| Meshing, Hole Filling/Smoothing, Data Merger | 48  |
| Stereo Views w/ All Enhancements             |     |
| Head-Tracked                                 | 36  |
| Head-Tracked w/Prediction (Render Rate)      | 34  |
| Head-Tracked w/Prediction (Barrier Rate)     | 74  |

## 6. Conclusions and Future Work

We have presented solutions to several issues related to building a 3D capture system using multiple depth cameras: resolving interference, data merging, and color matching between units. We have also introduced an eye position tracking system using depth sensors and demonstrated effective stereo display using rendering rates that would not usually support significant head motion. We have also incorporated improvements in calibration, data filtering, and data merger that served to improve image quality over a previous version of the system [12].

Using these solutions, we have demonstrated a telepresence system that is able to capture a fully dynamic 3D scene the size of a cubicle while allowing a remote user to look around the scene from any viewpoint. The system preserves eye gaze and does not require the user to wear any encumbrances. Using a single PC and graphics card, our system was able to render head-tracked stereo views at interactive rates and maintained stereo percept even with moderate head movement speeds.

Although our test system is functional, there are areas that we would like to improve, notably image quality. Although we incorporated a simple temporal noise suppression function in our system, some temporal noise artifacts are still present at the edges of objects where the depth camera alternates between providing a value and reporting no data at a given pixel. These depth pixels could be modified to keep a steady state or object contours could be smoothed and gaps could be filled in. Color calibration could be enhanced by combining our color correspondence-building algorithm with more robust color matching functions.

We also intend to expand our test setup into the “ideal” system shown in the bottom of Figure 3 by supporting 3D capture and 3D display for multiple users in both spaces. As seen in Figure 2, we already support 3D capture of multiple users. In this future system, we intend to add support for multiple tracked users on both sides.

Finally, we would like to expand the communication ability of our system by adding support for virtual objects that can be manipulated naturally by persons in the scene. Figure 21 shows an early experiment.

## Acknowledgments

The authors would like to thank Herman Towles for proposing this system, Andrei State for helping to produce the supplemental video and both for making helpful suggestions for improving this paper. We also thank Kurtis Keller for advice on reducing Kinect interference, Adrian Ilie for recommendations regarding color calibration, Ava Pope for helping to measure the wavelength of the Kinect’s IR lasers. This work was supported in part by the National Science Foundation (award CNS-0751187) and by the BeingThere Centre, a collaboration of UNC Chapel Hill, ETH Zurich, NTU Singapore, and the Media Development Authority of Singapore.



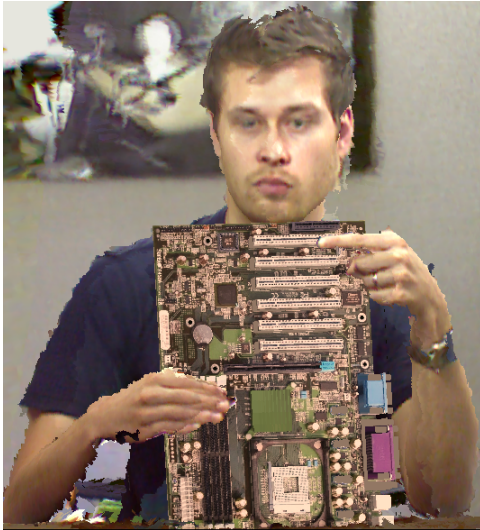


Figure 21: Mixed Reality Application. A 3D virtual object (circuit board) is incorporated into the scene during real-time 3D capture and naturally occludes and is occluded by real objects.

## References

- 976
- 977 [1] R. Raskar, G. Welch, M. Cutts, A. Lake, L. Stesin, H. Fuchs,  
978 The office of the future: a unified approach to image-based mod-  
979 eling and spatially immersive displays, in: Proceedings of the 25th  
980 annual conference on Computer graphics and interactive techniques,  
981 SIGGRAPH '98, ACM, New York, NY, USA, 1998, pp. 179–188.  
982 doi:<http://doi.acm.org/10.1145/280814.280861>.  
983 URL <http://doi.acm.org/10.1145/280814.280861>
- 984 [2] S. J. Gibbs, C. Arapis, C. J. Breiteneder, Teleport towards  
985 immersive copresence, *Multimedia Systems* 7 (1999) 214–221,  
986 10.1007/s005300050123.  
987 URL <http://dx.doi.org/10.1007/s005300050123>
- 988 [3] W. Matusik, H. Pfister, 3d tv: a scalable system for real-  
989 time acquisition, transmission, and autostereoscopic display  
990 of dynamic scenes, *ACM Trans. Graph.* 23 (2004) 814–824.  
991 doi:<http://doi.acm.org/10.1145/1015706.1015805>.  
992 URL <http://doi.acm.org/10.1145/1015706.1015805>
- 993 [4] O. Schreer, I. Feldmann, N. Atzpadin, P. Eisert, P. Kauff, H. Belt, 3d-  
994 presense -a system concept for multi-user and multi-party immersive 3d  
995 videoconferencing, in: *Visual Media Production (CVMP 2008)*, 5th Eu-  
996 ropean Conference on, 2008, pp. 1–8.
- 997 [5] A. Jones, M. Lang, G. Fyffe, X. Yu, J. Busch, I. McDowall, M. Bo-  
998 las, P. Debevec, Achieving eye contact in a one-to-many 3d video  
999 teleconferencing system, *ACM Trans. Graph.* 28 (2009) 64:1–64:8.  
1000 doi:<http://doi.acm.org/10.1145/1531326.1531370>.  
1001 URL <http://doi.acm.org/10.1145/1531326.1531370>
- 1002 [6] T. Balogh, P. T. Kovács, Real-time 3d light field transmission, Vol. 7724,  
1003 SPIE, 2010, p. 772406. doi:10.1117/12.854571.  
1004 URL <http://link.aip.org/link/?PSI/7724/772406/1>
- 1005 [7] C. Kuster, T. Popa, C. Zach, C. Gotsman, M. Gross, Freecam: A hybrid  
1006 camera system for interactive free-viewpoint video, in: *Proceedings of*  
1007 *Vision, Modeling, and Visualization (VMV)*, 2011.
- 1008 [8] M. Gross, S. Würmlin, M. Naef, E. Lamboray, C. Spagno, A. Kunz,  
1009 E. Koller-Meier, T. Svoboda, L. Van Gool, S. Lang, K. Strehle,  
1010 A. V. Moere, O. Staadt, blue-c: a spatially immersive display and 3d  
1011 video portal for telepresence, *ACM Trans. Graph.* 22 (2003) 819–827.  
1012 doi:<http://doi.acm.org/10.1145/882262.882350>.  
1013 URL <http://doi.acm.org/10.1145/882262.882350>
- 1014 [9] G. Kurillo, R. Bajcsy, K. Nahrsted, O. Kreylos, Immersive 3d envi-  
1015 ronment for remote collaboration and training of physical activities, in:  
1016 *Virtual Reality Conference, 2008. VR '08. IEEE*, 2008, pp. 269–270.  
1017 doi:10.1109/VR.2008.4480795.
- 1018 [10] H. H. Baker, N. Bhatti, D. Tanguay, I. Sobel, D. Gelb, M. E.  
1019 Goss, W. B. Culbertson, T. Malzbender, Understanding perfor-  
1020 mance in coliseum, an immersive videoconferencing system, *ACM*  
1021 *Trans. Multimedia Comput. Commun. Appl.* 1 (2005) 190–210.  
1022 doi:<http://doi.acm.org/10.1145/1062253.1062258>.  
1023 URL <http://doi.acm.org/10.1145/1062253.1062258>
- 1024 [11] B. Petit, J.-D. Lesage, C. Menier, J. Allard, J.-S. Franco, B. Raffin,  
1025 E. Boyer, F. Faure, Multicamera real-time 3d modeling for telepresence  
1026 and remote collaboration, *International Journal of Digital Multimedia*  
1027 *Broadcasting* 2010 (2009) 247108–12.
- 1028 [12] A. Maimone, H. Fuchs, Encumbrance-free telepresence system with real-  
1029 time 3d capture and display using commodity depth cameras, in: *Mixed*  
1030 *and Augmented Reality (ISMAR)*, 2011 10th IEEE International Sym-  
1031 posium on, 2011, pp. 137–146. doi:10.1109/ISMAR.2011.6092379.
- 1032 [13] L. Muhlbach, M. Bocker, A. Prussog, Telepresence in videoconferencing:  
1033 a study on stereoscopy and individual eye contact, *Human Factors*  
1034 37 (1995) 16.
- 1035 [14] B. Freedman, A. Shpunt, M. Machline, Y. Arieli, Depth mapping using  
1036 projected patterns, Patent Application, uS 2010/0118123 A1 (05 2010).
- 1037 [15] M. Meehan, Physiological reaction as an objective measure of presence in  
1038 virtual environments, Ph.D. thesis, University of North Carolina at Chapel  
1039 Hill (2001).
- 1040 [16] Z. Zhang, Flexible camera calibration by viewing a plane from unknown  
1041 orientations, in: *Computer Vision, 1999. The Proceedings of the Sev-*  
1042 *enth IEEE International Conference on*, Vol. 1, 1999, pp. 666–673 vol.1.  
1043 doi:10.1109/ICCV.1999.791289.
- 1044 [17] A. Maimone, H. Fuchs, A first look at a telepresence system with room-  
1045 sized real-time 3d capture and large tracked display, in: *Artificial Reality*  
1046 *and Telexistence (ICAT)*, The 21st International Conference on, 2011.
- 1047 [18] K. Berger, K. Ruhl, C. Brümmer, Y. Schröder, A. Scholz, M. Magnor,  
1048 Markerless motion capture using multiple color-depth sensors, in: *Proc.*  
1049 *Vision, Modeling and Visualization (VMV)*, 2011.
- 1050 [19] P. Merrell, A. Akbarzadeh, L. Wang, P. Mordohai, J.-M. Frahm, R. Yang,  
1051 D. Nister, M. Pollefeys, Real-time visibility-based fusion of depth maps,  
1052 in: *Computer Vision, 2007. ICCV 2007. IEEE 11th International Confer-*  
1053 *ence on*, 2007, pp. 1–8. doi:10.1109/ICCV.2007.4408984.
- 1054 [20] M. McGuire, A fast, small-radius gpu median filter, in: *ShaderX6*, 2008.  
1055 URL [http://graphics.cs.williams.edu/papers/](http://graphics.cs.williams.edu/papers/MedianShaderX6/)  
1056 [MedianShaderX6/](http://graphics.cs.williams.edu/papers/MedianShaderX6/)
- 1057 [21] C. Tomasi, R. Manduchi, Bilateral filtering for gray and color images, in:  
1058 *Computer Vision, 1998. Sixth International Conference on*, 1998, pp. 839  
1059 –846. doi:10.1109/ICCV.1998.710815.
- 1060 [22] H. Kawata, T. Kanai, Image-based point rendering for multiple range im-  
1061 ages, in: *Proceedings of the 2nd International Conference on Information*  
1062 *Technology for Application (ICITA 2004)*, Macquarie Scientific Publish-  
1063 ing, Sydney, 2004, pp. 478–483.
- 1064 [23] A. Ilie, G. Welch, Ensuring color consistency across multiple cameras,  
1065 in: *Proceedings of the Tenth IEEE International Conference on Computer*  
1066 *Vision - Volume 2, ICCV '05, IEEE Computer Society, Washington, DC,*  
1067 *USA, 2005*, pp. 1268–1275. doi:<http://dx.doi.org/10.1109/ICCV.2005.88>.  
1068 URL <http://dx.doi.org/10.1109/ICCV.2005.88>
- 1069 [24] M. A. Fischler, R. C. Bolles, Random sample consensus: a  
1070 paradigm for model fitting with applications to image analysis  
1071 and automated cartography, *Commun. ACM* 24 (1981) 381–395.  
1072 doi:<http://doi.acm.org/10.1145/358669.358692>.  
1073 URL <http://doi.acm.org/10.1145/358669.358692>
- 1074 [25] Q. Cai, D. Gallup, C. Zhang, Z. Zhang, 3d deformable face tracking  
1075 with a commodity depth camera, in: *Proceedings of the 11th European*  
1076 *conference on computer vision conference on Computer vision: Part III,*  
1077 *ECCV'10, Springer-Verlag, Berlin, Heidelberg, 2010*, pp. 229–242.  
1078 URL [http://portal.acm.org/citation.cfm?id=1927006.](http://portal.acm.org/citation.cfm?id=1927006.1927026)  
1079 [1927026](http://portal.acm.org/citation.cfm?id=1927006.1927026)
- 1080 [26] P. Viola, M. Jones, Rapid object detection using a boosted cascade of sim-  
1081 ple features, in: *Computer Vision and Pattern Recognition, 2001. CVPR*  
1082 *2001. Proceedings of the 2001 IEEE Computer Society Conference on*,  
1083 Vol. 1, 2001, pp. 1–511 – I–518 vol.1. doi:10.1109/CVPR.2001.990517.
- 1084 [27] R. Lienhart, J. Maydt, An extended set of haar-like features for rapid  
1085 object detection, in: *Image Processing, 2002. Proceedings. 2002 In-*  
1086 *ternational Conference on*, Vol. 1, 2002, pp. I–900 – I–903 vol.1.  
1087 doi:10.1109/ICIP.2002.1038171.
- 1088 [28] M. Castrillón Santana, O. Déniz Suárez, M. Hernández Tejera,  
1089 C. Guerra Artal, Encara2: Real-time detection of multiple faces at differ-

- 1090 ent resolutions in video streams, *Journal of Visual Communication and*  
1091 *Image Representation* (2007) 130–140.
- 1092 [29] G. Ye, A. State, H. Fuchs, A practical multi-viewer tabletop au-  
1093 tostereoscopic display, in: *Mixed and Augmented Reality (ISMAR)*,  
1094 2010 9th IEEE International Symposium on, 2010, pp. 147 –156.  
1095 doi:10.1109/ISMAR.2010.5643563.
- 1096 [30] T. Peterka, R. Kooima, J. Girado, J. Ge, D. Sandin, A. Johnson, J. Leigh,  
1097 J. Schulze, T. DeFanti, Dynallax: Solid state dynamic parallax barrier  
1098 autostereoscopic vr display, in: *Virtual Reality Conference, 2007. VR*  
1099 *'07. IEEE, 2007*, pp. 155 –162. doi:10.1109/VR.2007.352476.
- 1100 [31] J. Kopf, M. F. Cohen, D. Lischinski, M. Uyttendaele, Joint bilateral upsampling, *ACM Trans. Graph.* 26.  
1101 doi:http://doi.acm.org/10.1145/1276377.1276497.  
1102 URL <http://doi.acm.org/10.1145/1276377.1276497>
- 1104 [32] H. Towles, W.-C. Chen, R. Yang, S.-U. Kum, H. F. N. Kelshikar, J. Mul-  
1105 ligan, K. Daniilidis, H. Fuchs, C. C. Hill, N. K. J. Mulligan, L. Holden,  
1106 B. Zeleznik, A. Sadagic, J. Lanier, 3d tele-collaboration over internet2,  
1107 in: *International Workshop on Immersive Telepresence, Juan Les Pins,*  
1108 2002.