

000
001
002
003
004
005
006
007
008
009054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Hybrid 3D Reconstruction Using Photometric Stereo and Photogrammetry

Yi Cheng Zhu, Cheng Xu, Miaoqi Zhang

Abstract

3D reconstruction of objects is a fundamental task in computer vision, with applications in areas such as robotics, entertainment, and cultural heritage preservation. In this paper, we propose a hybrid approach that combines photogrammetry and near-light photometric stereo to achieve high-quality 3D reconstructions. Our method leverages a reference camera for acquiring images of the object under varying illumination conditions and a floating camera for capturing images from different angles. The pose of the floating camera is estimated using the Scale-Invariant Feature Transform (SIFT), Fast Library for Approximate Nearest Neighbors (FLANN), and Random Sample Consensus (RANSAC) algorithms. The 3D reconstruction is then achieved by integrating the surface normal map generated using the near-light photometric stereo approach and the point cloud obtained from the photogrammetry pipeline. Our experiments demonstrate the feasibility and potential of such a method.

1. Introduction

3D reconstruction of objects has been a topic of interest in computer vision for many years. Traditional methods for 3D reconstruction include photogrammetry, which uses multiple overlapping images to estimate the structure and geometry of an object, and photometric stereo, which estimates surface normals from images taken under different illumination conditions. Both of these methods have their strengths and weaknesses; photogrammetry can provide accurate geometry but may struggle with textureless or reflective surfaces, while photometric stereo can recover fine surface details but is sensitive to lighting conditions and may not provide accurate depth information.

In this work, we propose a novel hybrid approach that combines the strengths of both photogrammetry and near-light photometric stereo to achieve high-quality 3D reconstructions. Our method employs a reference camera for acquiring images of the object and a floating camera for capturing images under varying illumination conditions. We estimate the pose of the floating camera and predict the point cloud of the scene at the same time with

the structure-from-motion algorithm. The photogrammetry pipeline focuses on using floating camera images for the point cloud generation, while the near-light photometric stereo approach generates a surface normal map from acquired data of the reference camera and floating camera's light source positions.

The resulting surface normal map and depth maps are then integrated to generate a complete 3D reconstruction of the object. This hybrid approach benefits from the combined strengths of both photometric stereo and photogrammetry, achieving higher quality reconstructions than either method alone. In this paper, we present the details of our approach and try to demonstrate its effectiveness through a series of experiments on real-world objects.

2. Related Works

It has been shown by our sister field of computer graphics that having both surface normal and point cloud data has numerous benefits, including better interpolation of the model [7] and greater robustness to perturbation in measurements. [2]

In the field of computer vision, however, we usually use one to calculate the other for ease of setup, either with Poisson Surface Reconstruction [5] or the reverse process. [6].

This project is intended to build upon the prior works of Yeh, Chia-Kai, et al. [10]. Where their team proposed the use of a stationary camera to both act as a grounded coordinate frame for obtaining the pose for photogrammetry and also the camera for photometric stereo. With this setup, and assuming that the point light source(flashlight) is close to the moving camera, the measurement process would be simplified to taking a few photos with the flash on. We plan on utilizing this process and exploring the potential of combining both surface normal and depth maps to obtain a more accurate measurement.

3. Theory

3.1. overall setup

The aim of the project is to create a workflow that demands as little input and setup as possible from the user. We have simplified the hardware requirement to two phones, a tripod, and a printed Aruco marker, and the workflow from

108 the user's perspective is as follows: the user will need to
109 set up one phone on the tripod facing downwards and start
110 recording the video. Then the user should set up the scene
111 with an Aruco marker in the frame. The user then only
112 needs to take 7-10 pictures with the other phone with the
113 flashlight on at various different angles. The pipeline should
114 then be able to extract all the information needed to create
115 a high-quality 3d reconstruction.
116

117 3.2. calibration

118 The calibration matrices(the camera intrinsic matrix and
119 the distortion coefficients) are needed at various points in
120 the calculations. Such matrices can usually be accessed
121 through APIs on majorities of Android and Apple smart-
122 phones, but in the event that they are not available, we also
123 built in a camera calibration function that requires the user
124 to take multiple pictures of ChAruCo boards at various an-
125 gles.
126

127 3.3. Photogrammetry

128 Photogrammetry is a technique that derives 3D infor-
129 mation from multiple overlapping images of an object or
130 scene. The underlying principles of photogrammetry are
131 rooted in projective geometry and the concept of structure-
132 from-motion (SfM). The key steps involved in a typical pho-
133 togrammetry pipeline include:
134

- 135 1. Camera calibration: Determining the intrinsic and ex-
136 trinsic parameters of the camera, such as the focal
137 length, distortion coefficients, and camera pose. Accu-
138 rate calibration is crucial for obtaining reliable 3D
139 reconstructions.
140
- 141 2. Feature detection and matching: Identifying key points
142 in the images and establishing correspondences be-
143 tween them to determine their relative positions and
144 orientations. This process usually involves the use
145 of robust feature detectors and descriptors, such as
146 SIFT or SURF, and efficient matching algorithms like
147 FLANN.
148
- 149 3. Camera pose estimation: The pose of a camera is de-
150 fined by its position and orientation in the 3D world.
151 Estimating the camera pose involves determining these
152 parameters for each image in the input sequence. This
153 is usually done by applying the essential matrix or the
154 fundamental matrix to estimate the relative pose be-
155 tween the camera pairs. The absolute pose of the cam-
156 eras is obtained by chaining the relative poses through
157 the sequence.
158
- 159 4. Triangulation: Once the camera poses are known, the
160 next step is to estimate the 3D position of the points
161 in the scene. This is done through triangulation, which

162 involves finding the intersection of the viewing rays
163 from multiple camera viewpoints for a given point cor-
164 respondence. The intersection of these rays gives the
165 3D position of the point in the world. Several methods
166 can be used for triangulation, such as linear triangu-
167 lation, direct linear transformation (DLT), or iterative
168 methods like Levenberg-Marquardt.
169

- 170 5. Bundle adjustment: The camera pose estimation and
171 triangulation steps might have errors due to noise in
172 the feature matching or inaccuracies in the pose es-
173 timation. Bundle adjustment is a global optimiza-
174 tion technique that refines the camera poses and 3D
175 points simultaneously to minimize the reprojection er-
176 ror, which is the difference between the observed 2D
177 positions of the points in the images and the positions
178 obtained by projecting the 3D points onto the images
179 using the estimated camera poses. This optimization
180 is typically performed using a non-linear least squares
181 method like Levenberg-Marquardt.
182
- 183 6. Dense reconstruction: Estimating depth information
184 for each pixel in the images by finding correspon-
185 dences between overlapping images. This step can be
186 achieved using multi-view stereo techniques, such as
187 semi-global matching (SGM), patch-based methods,
188 or depth map fusion algorithms.
189

190 3.4. Near-Light Photometric Stereo

191 The Near-Light Photometric Stereo approach is based on
192 the work of Woodham and further developed by other re-
193 searchers. This method estimates the surface normal map
194 of an object using acquired data from the reference camera
195 and the floating camera's light source positions. The key
196 steps involved in the near-light photometric stereo pipeline
197 are:
198

- 199 1. Illumination vector computation: For each floating
200 camera image, the light source position is approxi-
201 mated using the camera's location or a pre-calculated
202 adjustment factor. Accurate estimation of the illumi-
203 nation vectors is essential for reliable surface normal
204 estimation.
205
- 206 2. Image preprocessing: The reference camera images
207 are preprocessed to reduce noise, correct for illumi-
208 nation inconsistencies, and enhance the photometric
209 stereo process. This step may involve normalization,
210 filtering, or other techniques, depending on the quality
211 of the images.
212
- 213 3. Pixel-wise intensity measurements: For each pixel
214 in the reference camera images, the intensity mea-
215 surements are collected to form a vector $I_i = (I_{i1}, I_{i2}, \dots, I_{in})$, where n is the number of images and
216 I_{ij} is the intensity of the i -th pixel in the j -th image.
217

- 216 4. Surface normal estimation: Using the illumination
 217 vectors L_i and the intensity measurements I_i , the
 218 surface normal \mathbf{N} for each pixel can be estimated using a
 219 least-squares approach. The surface normals can then
 220 be combined to create a surface normal map for the
 221 object.
 222
- 223 5. Integration and fusion: The surface normal map and
 224 the depth maps obtained from the photogrammetry
 225 pipeline are integrated to generate a complete 3D re-
 226 construction of the object. This may involve a multi-
 227 view stereo approach or other fusion techniques to
 228 ensure consistency and accuracy in the final output.
 229 The resulting 3D model benefits from the combined
 230 strengths of both photometric stereo and photogram-
 231 metry, achieving higher quality reconstructions than
 232 either method alone.

233 4. Analysis

234 4.1. Photogrammetry Analysis

235 The photogrammetry process can be analyzed based on
 236 the principles of projective geometry and structure-from-
 237 motion (SfM). SfM is a technique for estimating the 3D
 238 structure of a scene and the camera motion using a set of
 239 2D images. In the SfM approach, camera poses and 3D
 240 points are estimated simultaneously, making it a non-linear
 241 optimization problem. This problem can be solved iter-
 242 atively using bundle adjustment, which refines the camera
 243 poses and 3D points by minimizing the reprojection error
 244 between observed and predicted image features.

245 4.1.1 Feature extraction and matching

246 The first step in SfM is to extract distinctive features (points,
 247 corners, or edges) from each image in the sequence. Com-
 248 mon feature extraction algorithms are the Scale-Invariant
 249 Feature Transform (SIFT), Oriented FAST and Rotated
 250 BRIEF (ORB), or Speeded-Up Robust Features (SURF).
 251 Once features are extracted from each image, we need to
 252 find correspondences between the features across the
 253 images. This is done using feature descriptors, which are
 254 vectors that describe a feature's local neighborhood. Fea-
 255 ture matching can be performed using methods like Nearest
 256 Neighbor Search or Lowe's ratio test.

257 4.1.2 Camera motion estimation

258 We can estimate the camera motion (rotation and transla-
 259 tion) between two images using the Essential Matrix (E) or
 260 the Fundamental Matrix (F). The Essential Matrix relates
 261 the corresponding points in two images when the camera's
 262 intrinsic parameters are known, while the Fundamental
 263 Matrix is used when the camera's intrinsic parameters are
 264 unknown.

265 For the Essential Matrix, let's denote two corresponding
 266 points in the images as p_1 and p_2 . We have:

$$267 p_2^T * E * p_1 = 0$$

268 where p_1 and p_2 are homogeneous coordinates of the
 269 corresponding points, and E is the Essential Matrix.

270 The Fundamental Matrix, F , is defined similarly:

$$271 p_2^T * F * p_1 = 0$$

272 If the camera's intrinsic parameters are known, we can
 273 compute the Essential Matrix (E) from the Fundamental
 274 Matrix (F) as:

$$275 E = K^T * F * K$$

276 where K is the camera's intrinsic matrix.

277 The Essential Matrix can then be decomposed into the
 278 camera's rotation matrix (R) and translation vector (t) using
 279 the Singular Value Decomposition (SVD) method.

280 4.1.3 Triangulation

281 Once the camera motion is estimated, we can find the 3D
 282 position of the matched feature points in the scene using
 283 triangulation. Given the camera matrices P_1 and P_2 for the
 284 two views, and their corresponding points x_1 and x_2 , we
 285 want to find the 3D point X that minimizes the reprojection
 286 error:

$$287 \text{argmin} \|x_1 - P_1 * X\|^2 + \|x_2 - P_2 * X\|^2$$

288 The solution can be found using the Direct Linear Trans-
 289 formation (DLT) algorithm or the Least Squares Triangula-
 290 tion method.

291 4.1.4 Bundle adjustment

292 After obtaining initial estimates of the camera motion and
 293 3D points, bundle adjustment is used to refine these esti-
 294 mates by minimizing the reprojection error across all im-
 295 ages and points. The reprojection error is defined as the
 296 difference between the observed image point and the pro-
 297 jection of the corresponding 3D point onto the image plane:

$$298 E = \text{sum}(\|x_{ij} - P_i * X_j\|^2)$$

299 where x_{ij} is the observation of the j -th point in the i -th
 300 image, P_i is the projection matrix of the i -th camera, and X_j
 301 is the 3D position of the j -th point. This is a non-linear op-
 302 timization problem and can be solved using the Levenberg-
 303 Marquardt algorithm.

304 The accuracy of the photogrammetry pipeline depends
 305 on various factors, such as the quality and resolution of the

324 images, the baseline (i.e., the distance between cameras),
 325 the number of images, and the quality of the feature matching
 326 and dense reconstruction algorithms used. By analyzing
 327 the performance of the photogrammetry process on different
 328 datasets and scenarios, one can identify the strengths
 329 and limitations of the method and explore potential im-
 330 provements.
 331

332 4.2. Near-Light Photometric Stereo Analysis

333 The near-light photometric stereo approach can be par-
 334 ticularly advantageous when dealing with objects that have
 335 textureless or reflective surfaces, as these characteristics
 336 can pose challenges for traditional photogrammetry meth-
 337 ods. By incorporating the near-light photometric stereo
 338 technique into the 3D reconstruction pipeline, the resulting
 339 model can achieve improved surface details, accurate geo-
 340 metry, and robust performance on a wider range of object
 341 types and surface properties.
 342

343 The Near-Light Photometric Stereo approach is based on
 344 the work of Woodham [9] and further developed by other
 345 researchers [3, 4]. This section outlines the process of gen-
 346 erating a surface normal map using the acquired data from
 347 the reference camera and the floating camera’s light source
 348 positions.
 349

350 4.2.1 Illumination Vectors

351 For each floating camera image, the light source position
 352 is approximated using the camera’s location or a pre-
 353 calculated adjustment factor. This results in a set of illumi-
 354 nation vectors, $L_i = (l_{ix}, l_{iy}, l_{iz})$, representing the direc-
 355 tion and intensity of the light source.
 356

357 4.2.2 Image Preprocessing

358 The reference camera images corresponding to each float-
 359 ing camera image will be preprocessed to reduce noise and
 360 enhance the photometric stereo process. This may involve
 361 normalization, filtering, or other techniques, depending on
 362 the quality of the images.
 363

364 4.2.3 Pixel-wise Intensity Measurements

365 For each pixel in the reference camera images, the inten-
 366 sity measurements are collected to form a vector $I_i =$
 367 $(I_{i1}, I_{i2}, \dots, I_{in})$, where n is the number of images and I_{ij}
 368 is the intensity of the i -th pixel in the j -th image.
 369

370 4.2.4 Surface Normal Estimation

371 We estimate the surface normals from photometric stereo
 372 following the methods proposed by [8].
 373

374 We first re-scale the world coordinates to pixel coordi-
 375 nates by multiplying by the constant
 376

$$377 d = \frac{\mu \hat{z}}{f} \\ 378$$

379 where f is the focal length, μ is the camera pixel size,
 380 and \hat{z} the approximate mean depth of the object from the
 381 camera.
 382

383 Let a single light source be at e_k , we describe a single
 384 light pencil from light source k at pixel p as
 385

$$386 \hat{L}_{pk} = \frac{L_k - X_p}{\|L_k - X_p\|^2} e_k \\ 387$$

388 where X_p is the 3D location of the pixel. The depth com-
 389 ponent of X_p is initially set to 1, and we iteratively solve
 390 the depth map.
 391

392 Using the illumination vectors L_i and the intensity mea-
 393 surements I_i , the surface normal \mathbf{N} for each pixel can be
 394 estimated using a least-squares approach:
 395

396 The image intensity at a pixel is then described as
 397

$$398 I_{pk} = B_p^T \hat{L}_{pk} \\ 399$$

400 where B_p is the surface normal at the location.
 401

402 We solve for B_p by computing the pseudo-inverse of
 403 \hat{L}_{pk} .
 404

405 Using B_p , we construct the gradient field by dividing out
 406 the Z component B_p .
 407

408 This is often a non-integrable gradient field.
 409

410 Following the methods of [1], we attempt to solve a plau-
 411 sible solution to this non-integrable gradient field.
 412

413 Let the two components of the gradient field be (p, q) .
 414 At each pixel we construct the 2×2 matrix H by gaussian
 415 blurring component-wise the matrix $\begin{bmatrix} p^2 & p \times q \\ p \times q & q^2 \end{bmatrix}$. Let
 416 u_1, u_2 be the eigenvalues of H , construct new eigenvalues as
 417

$$418 \lambda_1 = \begin{cases} 1 & u_1 = 0 \\ \beta + 1 - \exp -3.315/u_1^4 & u_1 > 0 \end{cases} \\ 419 \lambda_2 = 1 \\ 420$$

421 Construct the matrix D using eigenvectors and eigenval-
 422 ues λ_1, λ_2 .
 423

424 Then let
 425

$$426 D = \begin{bmatrix} d_{11} & d_{12} \\ d_{21} & d_{22} \end{bmatrix} \\ 427$$

$$428 u_D = \text{div}(d_{11}p + d_{12}q, d_{21}p + d_{22}q) \\ 429$$

430 We solve for the new depth map Z as
 431

$$432 \nabla_D^2 Z = u_D \\ 433$$

434 Then, by iteratively updating the normal map, then inte-
 435 grating the resulting gradient field for a new depth map, it
 436 converges into the accurate final depth map.
 437

432

4.2.5 Integration and Fusion

433

The surface normal map and the depth maps obtained from the photogrammetry pipeline will be integrated to generate a complete 3D reconstruction of the object. This may involve a multi-view stereo approach or other fusion techniques to ensure consistency and accuracy in the final output. The resulting 3D model will benefit from the combined strengths of both photometric stereo and photogrammetry, achieving higher quality reconstructions than either method alone.

443

The accuracy of the near-light photometric stereo approach is influenced by several factors, including the quality and consistency of the illumination vectors, the precision of the surface normal estimation, and the robustness of the integration and fusion techniques. By examining the performance of the method on various datasets and under different conditions, one can assess its capabilities and identify areas for improvement.

451

5. experimental results

452

5.1. Frame extraction

453

In order for the hybrid method to function, we need to establish a mapping between the pictures from the overhead camera and the portable camera. We have decided to record the overhead camera as a video so as to not worry about the temporal syncing of the shutter until the preprocessing step. At the preprocessing step, we first evaluate the light level at each frame (see figure. 1). Then we extract the highest n amount of frames with the restrain that the frames will not be within 0.2 seconds of each other and n being the number of photos the portable camera took. We then convolve the neighboring two frames with the selected frame to account for the rolling shutter effect to finalize our preprocessing step. The corresponding pairs of images can be seen in figure. 8.

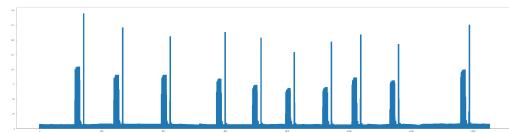


Figure 1. light level over each frame of the video

477

5.2. calibration

480

To provide a tool for calibration for cases where the internal matrix and the distortion coefficients are not accessible, we created a supplementary process to calibrate the cameras. We used the existing implementation of charuco boards in opencv to perform the calibration. The result and validations can be seen in figure. 2.

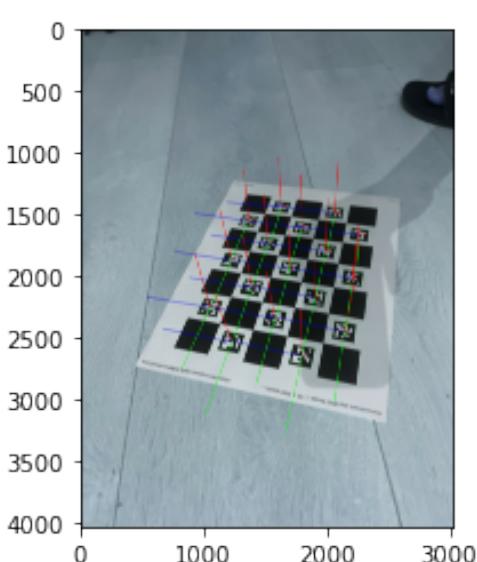


Figure 2. calibration done with charuco boards

5.3. pose estimation

Although the poses can be obtained by Structure from Motion, to increase robustness, we provided to option to use an aruco board as the poses for the photometric stereo step. We chose a 2x2 aruco board to ensure a more robust measurement when compared to a single marker. We once again leverage the libraries provided by opencv for this function. The unit test for the submodule can be seen in figure. 3.

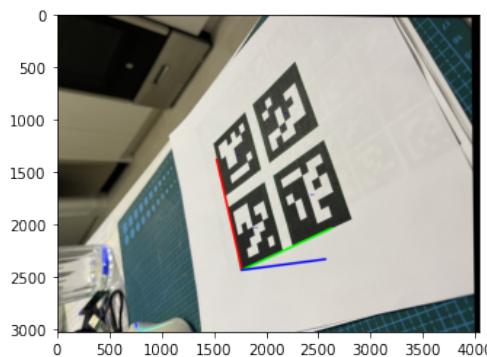


Figure 3. pose estimation using aruco board

5.4. photogrammetry

Our photogrammetry pipeline managed to extract the correct point cloud for the scene in the dataset(see figure. 4). Unfortunately we are unable to implement the dense estimation step in time. as a result, our depth map is not able to get per-pixel depth value. We will touch on that in the depth estimation section.

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

A 3D point cloud visualization showing a dense cluster of points forming a complex, irregular shape, likely representing a reconstructed surface or volume.

Figure 4. the point cloud output of our photogrammetry algorithm

5.5. depth estimation

As previously mentioned, our point cloud is not dense enough to directly convert to a depth map. As a result, we had to linearly interpolate the depth map for pixels where there are sufficient neighboring data. The result of this step can be seen in figure. 5.

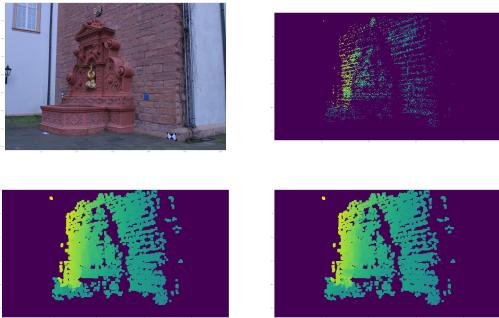


Figure 5. the depth map directly from projecting and after interpolation.

5.6. photometric stereo

The result for photometric stereo can be seen in Figure 6. and Figure 7.

6. concluding remarks

We have shown throughout this project that it is plausible to create a pipeline that incorporates both photogrammetry and photometric stereo while still maintaining the simplicity of only photogrammetry. Unfortunately, we were unable to implement the dense point cloud estimation part of photogrammetry, which created a rift in the entire pipeline and prevented us from testing it end to end. We were also unable to correlate the scale of the scene with its real-world measurements without the help of aruco markers. Those are the

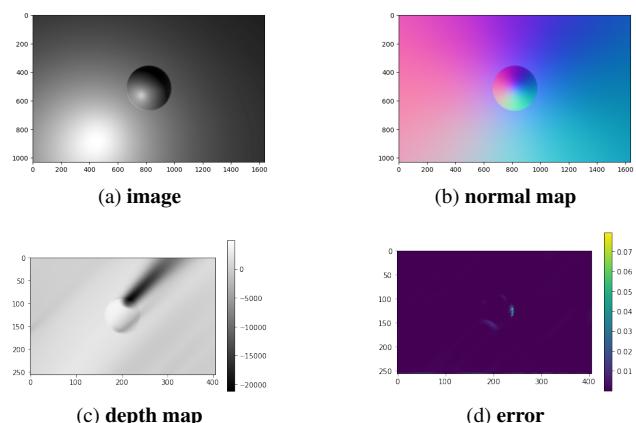


Figure 6. The result from photometric stereo

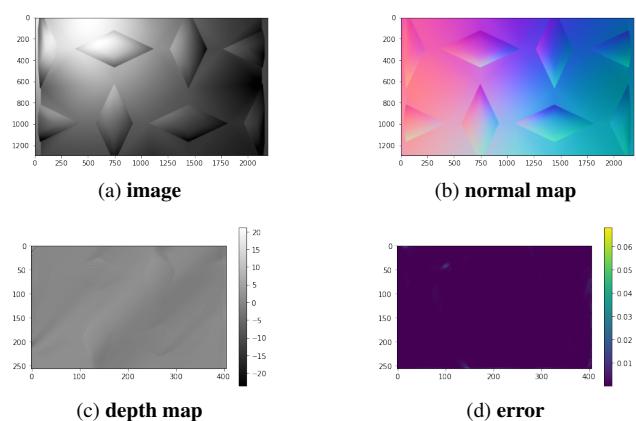


Figure 7. The result from photometric stereo

areas that warrant further exploration.

648

References

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

