

VOLUMETRIC RECONSTRUCTION APPLIED TO PERCEPTUAL STUDIES OF SIZE AND WEIGHT

J. BALZER, M. PETERS, AND S. SOATTO

ABSTRACT. We explore the application of volumetric reconstruction from structured-light sensors in cognitive neuroscience, specifically in the quantification of the *size-weight illusion*, whereby humans tend to systematically perceive smaller objects as heavier. We investigate the performance of two commercial structured-light scanning systems in comparison to one we developed specifically for this application. Our method has two main distinct features: First, it only samples a sparse series of viewpoints, unlike other systems such as the Kinect Fusion. Second, instead of building a distance field for the purpose of points-to-surface conversion directly, we pursue a first-order approach: the distance function is recovered from its gradient by a *screened Poisson reconstruction*, which is very resilient to noise and yet preserves high-frequency signal components. Our experiments show that the quality of metric reconstruction from structured light sensors is subject to systematic biases, and highlights the factors that influence it. Our main performance index rates estimates of volume (a proxy of size), for which we review a well-known formula applicable to incomplete meshes. Our code and data will be made publicly available upon completion of the anonymous review process.

1. INTRODUCTION

1.1. **Motivation.** We like to believe that our sensory systems provide us with precise and accurate information about objects within the environment, but our perception is often subject to systematic errors, or *illusions* (Figure 1). These can also occur between sensory modalities, often with visual information influencing haptic (touch) estimates of properties such as size or weight. For example, a curious experience occurs when we lift two objects of equal weight but different size; systematically and repeatably, the smaller object feels heavier than the larger. This *size-weight illusion* (SWI) [3] cannot be explained by simple motor force error (i.e., it is not simply due to the production of more lifting or grip force for the larger object) [6, 9], and so carries important implications for the dynamics of sensory integration between vision and haptics. Likewise, altered visual appearance of an object (e.g. through stereoscopic goggles [4] or optical distortion with prisms [18]) can significantly impact haptically-judged estimates of its size. Simply put, when an object looks bigger than it really is, it feels bigger, too – and any mismatch between vision and touch often goes completely unnoticed.

In order to establish a solid quantitative empirical assessment of these illusions, we have developed methodologies to examine the relationship between true size and perceived size. Previous investigations have uncovered evidence that the relationship between an object’s true volume and its perceived volume often follows a power function with an average exponent of 0.704 ($\sigma = 0.08$) [7]. However, these prior investigations have predominantly used objects which are geometric, symmetrical,

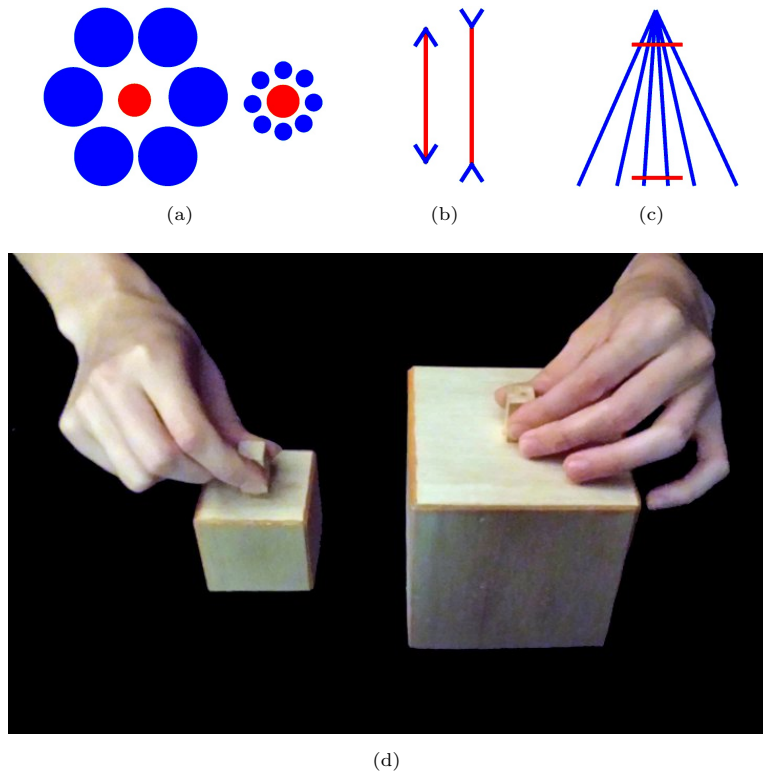


FIGURE 1. The (a) Ebbinghaus, (b) Müller-Lyer, and (c) Ponzo illusions demonstrate that perception is not always veridical [8]. For each, the orange elements are identical, although they appear unequal. (d) In the *size-weight illusion*, two objects that weigh the same feel as if the smaller one is heavier [3].

and convex – properties which alone cannot adequately capture the range of objects regularly encountered in everyday environments. Thus, to systematically and comprehensively explore the relationship between true volume and perceived volume so as to better understand this percept’s contribution to visual-haptic integration and, consequently, perception in general, we have developed dedicated methods to capture an ecologically valid set of stimuli.

Our goal is to build a data base of digital models of our specimens, which would allow us to infer volume and any other desired geometric properties. A method to create such models should meet the following list of criteria:

- It is mandatory that the sensing modality of choice be contactless: Cell phones and other hand-held consumer electronics, e.g., are among the classes of objects relevant to our psychological studies. They cannot simply be sunk in a fluid leveraging Archimedes’ principle, neither can anything which is permeated by the fluid, as doing so would not provide the desired *visible* volume.



FIGURE 2. We compare the performance of two scanning systems with the one proposed here at hand of a set of geometric primitives, whose volume can be easily measured by hand. Although shown as a collection, all items above have been reconstructed and texture-mapped individually.

- The underlying sensor should be inexpensive and easy to use for non-experts; e.g., researchers in psychology and neuroscience. This rules out dedicated lab equipment, e.g., for white-light interferometry and the like.
- The resulting models should exhibit some topological structure: For volume computations, it must at least be possible to distinguish interior and exterior. Therefore, point clouds alone are insufficient in this regard.
- Given the complexity of everyday objects – and that they frequently depart from the cubic, spherical, or cylindrical – we target an improvement of accuracy over back-of-the-envelope estimates.
- A rich spectrum of object classes is covered in terms of admissible geometry and reflectance properties. Specular surfaces, e.g., would need to be coated with powder to make them amenable to laser scanning. But this would contravene the first criterion and thus eliminates laser scanning from the list of candidates.

In light of these requirements, we opt for triangulation based on structured-light encoding as the primary sensing modality. The principle behind this method has been known for decades but has seen a renaissance in computer vision ever since Prime-sense introduced a fully functional color-range (“RGBD”) sensor unit in integrated-circuit design. This system-on-chip later became the core component of Microsoft’s Kinect, which subsequently had an considerable impact in geometry reconstruction, tracking, occlusion detection, and action recognition, among other applications.

1.2. Contribution and overview. We develop a system for structured-light scanning of small- to medium-scale objects, which we dub *Yet Another Scanner*, or YAS. Naturally, the question arises why we would need yet another scanner when several systems and commercial products are already available, e.g., Kinect Fusion [17],

ReconstructMe¹, Artec Studio², KScan3D³, Scanect⁴, Scenect⁵, and Fablitec’s 3d scanner⁶; in particular, when most of these generate visually highly-pleasing results.

The main reason is that, albeit visually pleasing, the reconstructions provided by these methods are subject to biases that make them unsuitable for scientific investigation. The analysis, which is presented in Sect. 3.2, compares the performance of YAS with that of two competing state-of-the-art implementations. While the reconstruction algorithm described in Sects. 2.2 and 2.3 itself is not novel, we carefully justify all choices to be made in its design w.r.t. above-listed requirements. Additionally, we address the issue that an aligned series of range images suffers from incompleteness precisely where the ground plane supports the object. Sect. 2.4 proposes a strategy to circumvent this problem in volume estimation which avoids complicated hole-filling algorithms. We believe that a tool which outperforms commercial software but is accessible for further scientific development may be of interest to the computer vision community as well. Hence, as the final contribution, we will distribute the source through the repository at <https://bitbucket.org/jbalzer/yas>.

2. SYSTEM DESCRIPTION

2.1. Data acquisition and calibration. In all our experiments studies, we used Microsoft’s Kinect and Primesense’s Carmine 1.09. Both devices are shipped with a factory calibration of depth and RGB camera intrinsics as well as the coordinate transformation between their local reference frames. Initial visual assessment (by the naked human eye) approves of the default calibration simply because the point clouds computed from the range image seem to be accurately colored by the values of the RGB image. Extensive tests, however, have shown that – in the spirit of our introductory remarks – such an evaluation is misleading, and significant metric improvements through manual re-calibration are possible. For this purpose, we rely on the toolbox accompanying the paper [11] to estimate all aforementioned parameters plus a depth uncertainty pattern, which varies both spatially and with depth itself.

2.2. View alignment. A calibrated sensor immediately delivers physically plausible depth data. The integration of measurements from different vantage points into a common 3-d model can thus be seen as the core challenge here. A comprehensive overview of the state of the art in scan alignment is found in the recent survey [19]. Essentially, one can distinguish between two approaches: *tracking* and *wide-baseline matching*. The former is at the heart of Kinect Fusion [17] and the majority of commercially available software. Its main motivation stems from the fact that correspondence is easier to establish when two images haven been acquired closely in time – supposing, of course, that the motion the camera has undergone between each image acquisition and the next meets certain continuity constraints. We believe, however, that for the purpose of small-scale object reconstruction, the disadvantages of tracking predominate. First and foremost, there is the question

¹<http://reconstructme.net/>

²<http://www.artec3d.com/software/>

³<http://www.kscan3d.com/>

⁴<http://skanect.manctl.com/>

⁵<http://www.faro.com/scenect/>

⁶<http://www.fablitec.com/>

of redundancy: How do we deal with the stream of depth data when operating an RGBD camera at frame rates up to 30 fps? On the one hand, redundancy is desirable because single depth images may not cover the entire surface of the unknown object, e.g., due to occlusions or radiometric disturbances of the projected infrared pattern. On the other hand, integration of overcomplete range data into a common 3-d model puts high demands on the quality of alignment. Most feature trackers operate on a reduced motion model⁷ and are thus prone to drift. Such deviations in combination with the uncertainty in the raw depth data can lead to a *stratification* of points in regions appearing in more than a single image. This effect is illustrated in Fig. 3(a), which becomes more severe with higher numbers of processed images.

Kinect Fusion [17] deals with redundancy by instantaneously merging the depth stream into an implicit surface representation over a probabilistic voxel grid. The extra dimension, however, raises memory consumption – even in implementations utilizing truncation or an efficient data structure such as an octree. Also, without a-priori knowledge of the specimen’s size, it is difficult to gauge the interplay between the spatial resolutions of 3-d grid and raw depth data, the latter being left exploited only sub-optimally. Last but not least, a system based on tracking is little user-friendly: It requires the operator to move the sensor as steadily as possible. Otherwise, temporal under-sampling or motion blur can lead to a total breakdown of the alignment process.

Here, we closely follow the wide-baseline matching procedure developed by the robotics community, notably in the work of Henry et al. [10]. It follows two quasi canonical steps: First, a set of local descriptors around interest points in each RGBD image is computed as well as a set of tentative matches between them. Such descriptors incorporate radiance and depth information either exclusively or in combination. Second, subsets of cardinality three are selected from all matches at random. Each subset admits a hypothesis about the rigid motion that transforms one of the point clouds into the other. The winning hypothesis is taken to be the one supporting the most matches, i.e., generating the highest number of *inliers* [5]. We review these initial two steps in the following sections.

2.2.1. *Sampling.* Let us formally consider the case of two views $l = 0, 1$, i.e., we look for two rigid motions $g_l \in \text{SE}(3)$, $g_l : \mathbf{x} \mapsto \mathbf{R}_l \mathbf{x} + \mathbf{t}_l$, with $\mathbf{R}_l \in \text{SO}(3)$ and $\mathbf{t}_l \in \mathbb{R}^3$. Without loss of generality, one can assume that g_0 coincides with the world reference frame, i.e., $\mathbf{R}_0 = \mathbf{I}$ and $\mathbf{t}_0 = \mathbf{0}$, which leads to a simpler notation of the unknowns $\mathbf{R}_1 = \mathbf{R}$ and $\mathbf{t}_1 = \mathbf{t}$.

A number of interest points $\{\mathbf{p}_0^i\}, \{\mathbf{p}_1^j\}$ with high response is extracted from each of the two RGB images corresponding to $g_{0,1}$ by means of the SIFT detector. The points are combined into a set of putative correspondences $\mathcal{C} = \{(\mathbf{p}_0^{i_k}, \mathbf{p}_1^{j_k}) \mid k = 1 \dots, n \in \mathbb{N}\}$ by thresholded forward-backward comparison of the distances between associated SIFT descriptors [15]. The search for nearest neighbors can be sped up by a locality-sensitive hashing technique or similar. However, we found in all of our experiments that the time consumed by a brute-force search was within acceptable limits. Next, we repeatedly draw a sample of three matches from \mathcal{C} and obtain a set of triples $\mathcal{H} = \{(k_1, k_2, k_3) \mid 1 \leq k_1, k_2, k_3 \leq n, k_1 \neq k_2 \neq k_3\}$.

2.2.2. *Consensus.* Implicitly, each of the elements of \mathcal{H} determines a hypothesis about the transformation we are looking for: Suppose we already know $g_1 \in \text{SE}(3)$,

⁷E.g., the Lucas-Kanade method assumes pure translational motion.

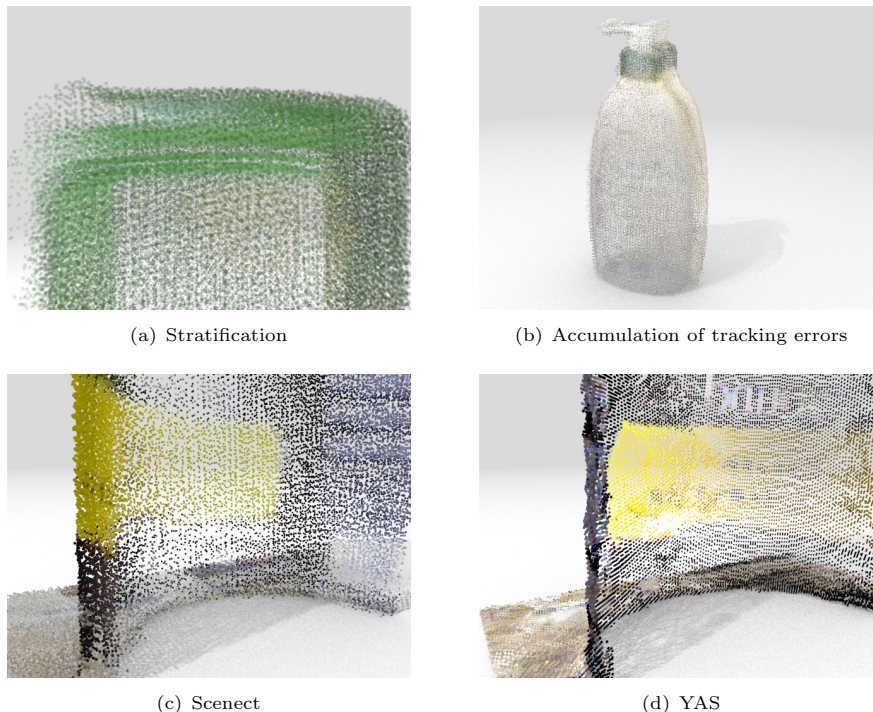


FIGURE 3. We promote a sparse sampling of viewpoints for covering the object of interest with depth map measurements, for (a) this reduces the impact of stratification and (b) reduces the probability of drift or tracking failure. (c)-(d) The difference in alignment quality only becomes noticeable in cross-sections of the point clouds.

then the geometric least-squares error for some $(k_1, k_2, k_3) \in \mathcal{H}$ is given by

$$e(k_1, k_2, k_3) = \sum_{k \in \{k_1, k_2, k_3\}} \frac{1}{2} \|\mathbf{x}_0^{i_k} - \mathbf{R}_k \mathbf{x}_1^{j_k} - \mathbf{t}_k\|^2. \quad (1)$$

Here, the six points $\mathbf{x}_0^{i_k}, \mathbf{x}_1^{j_k} \in \mathbb{R}^3$ equal the backprojections of the three matches $(\mathbf{p}_0^{i_k}, \mathbf{p}_1^{j_k})$ forming the current hypothesis. Given the intrinsic camera parameters, they can be easily computed from the data delivered by the calibrated depth sensor. Conversely, given a triple $(k_1, k_2, k_3) \in \mathcal{H}$, we can find a global minimizer g_1^* of the convex function (1) in the following way: Denote by $\bar{\mathbf{x}}_0$ the mean of $\mathbf{x}_0^{i_k}$ over k and define $\bar{\mathbf{x}}_0^{i_k} = \mathbf{x}_0^{i_k} - \bar{\mathbf{x}}_0$. The quantities $\bar{\mathbf{x}}_1$ and $\bar{\mathbf{x}}_1^{j_k}$ are defined analogously. A minimizer in the *entire general linear group* of matrices $\text{GL}(3)$ is

$$\mathbf{H} = \sum_{k \in \{k_1, k_2, k_3\}} \bar{\mathbf{x}}_0^{i_k} (\bar{\mathbf{x}}_1^{j_k})^\top,$$

cf. [20]. One needs to make sure that the optimal g_1^* involves a genuine *rotation matrix* by projecting \mathbf{H} onto $\text{SO}(3)$. This is commonly achieved by *Procrustes analysis*, essentially consisting of a singular-value decomposition: Write \mathbf{H} as the product $\mathbf{U}\Sigma\mathbf{V}^\top$ with two orthogonal factors $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{3 \times 3}$, then $\mathbf{R}_k^* = \mathbf{V}\mathbf{U}^\top$. Once

\mathbf{R}_k^* is known, the optimal translation vector can be computed as $\mathbf{t}_k^* = \bar{\mathbf{x}}_0 - \mathbf{R}_k^* \bar{\mathbf{x}}_1$. The transformation for the element of \mathcal{H} attaining the highest consensus among *all* matches in \mathcal{C} constitutes the solution to the global alignment problem [5]. The result is refined based on the inlier correspondences with the iterative closest-point (ICP) method [1]. This also ensures a *geometrically* continuous alignment, which is not guaranteed because \mathcal{H} was generated merely based on *photometry*.

2.3. Surface reconstruction. The common point cloud obtained after merging all aligned depth maps carries no information about the topological relationship between its elements, but as we will see shortly, such information plays a crucial part in volume estimation. There exists a wealth of algorithms for point-to-surface conversion, most of which depend on the signed or unsigned Euclidean distance field $\varphi : \mathbb{R}^3 \rightarrow \mathbb{R}$ induced by the point cloud. Kinect Fusion, e.g., computes φ directly. Alternatively, when the point cloud is oriented, i.e., each point \mathbf{x} is endowed with an estimate of the normal vector \mathbf{n} the surface should have at that location, one can search for the function φ minimizing

$$\int_D \frac{1}{2} \|\nabla\varphi - \mathbf{n}\|^2 d\mathbf{x} \rightarrow \min. \quad (2)$$

Here, with slight abuse of notation, \mathbf{n} refers to an (arbitrary) continuation of the normal field from the point set to some sufficiently large rectangular domain $D \subset \mathbb{R}^3$. Since the gradient of any scalar function is orthogonal to its level sets, this gives a family of *integral surfaces*

$$\Gamma = \{\mathbf{x} \in \mathbb{R}^3 \mid \varphi(\mathbf{x}) = C\} \quad (3)$$

of \mathbf{n} . A minimizer of (4) is found as the solution of the Euler-Lagrange equation

$$\Delta\varphi = \operatorname{div} \mathbf{n} \quad (4)$$

under natural boundary conditions (here of Neumann type). Eq. (4) is the well-known Poisson equation and eponymous for the *Poisson reconstruction* algorithm proposed in [12].

The motivation for increasing the order of differentiation as compared to the direct approach (i.e., that followed in Kinect Fusion) is twofold: First, Eq. (2) is a variant of the Dirichlet energy, which implies that small holes in the point cloud, i.e., areas where $\mathbf{n} = \mathbf{0}$, will automatically be in-painted harmonically. Second, for a solution to exist in the strong sense $\nabla\varphi = \mathbf{n}$, the normal field must be *integrable* or *curl-free*. Noise, which is very common in RGBD images, is responsible for most of the non-integrability in a measured normal field \mathbf{n} . In the variational setting (2), however, \mathbf{n} is implicitly replaced by the next-best gradient (its so-called *Hodge projection*, cf. [2]), which makes the approach very resilient to stochastic disturbances but at the same time destroys fine details. The smoothing effect can be mitigated by imposing Dirichlet conditions on (4) at a sparse set of salient points. This so-called *screened Poisson reconstruction* has recently been introduced in [13].

The point cloud is easily oriented exploiting the known topological structure of the pixel lattice: Given a depth parametrization of the surface $z(x, y)$ over the two orthogonal camera coordinate directions x and y , the normal can be written as $\mathbf{n} = (-\partial_x z, -\partial_y z, 1)^\top$. The partial derivatives of z w.r.t. camera and image

coordinates $(x, y)^\top$ respectively $(u, v)^\top$ are related by the chain rule:

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial u} \frac{\partial u}{\partial x} = \frac{f_u}{z} \frac{\partial z}{\partial u}, \quad \frac{\partial z}{\partial y} = \frac{\partial z}{\partial v} \frac{\partial v}{\partial y} = \frac{f_v}{z} \frac{\partial z}{\partial v}. \quad (5)$$

Here, f_u, f_v are the focal lengths of the pinhole depth camera, and finite-differencing provides an approximation to the gradient of $z(u, v)$. For the details of numerically solving (4) and selecting the constant C in (3) appropriately, we refer the reader to the original paper [12].

2.4. Volume estimation.

2.4.1. *Closed surfaces.* Suppose for the moment that Γ given by (3) is compact and closed. The volume V of the domain $\Omega \subset \mathbb{R}^3$ it encompasses is defined as the integral of the characteristic function χ_Ω of Ω :

$$V = |\Omega| = \int_{\mathbb{R}^3} \chi_\Omega d\Omega = \int_{\Omega} 1 d\Omega. \quad (6)$$

Unfortunately, an evaluation of this integral is not very practical for two reasons. First, doing so would require a regular grid over Ω , which introduces undesirable artefacts where it interacts with a discrete version of Γ . Second, an expensive nearest-neighbor problem would need to be solved⁸ to determine whether a point is inside or outside of Ω . The following trick is based on the classic Gauss divergence theorem, cf. [16], which relates the flow of *any* continuously differentiable vector field $\mathbf{v} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ through the boundary $\Gamma = \partial\Omega$ of Ω with its source density or *divergence* in the interior:

$$\int_{\Omega} \operatorname{div} \mathbf{v} d\Omega = \int_{\Gamma} \langle \mathbf{v}, \mathbf{n} \rangle d\Gamma. \quad (7)$$

The left-hand side of this equation does not quite resemble the right-hand side of (6), yet. However, this can be achieved by a clever choice of \mathbf{v} , e.g., $\mathbf{v} := (x, 0, 0)^\top$, but note that several variants will work equally well and that \mathbf{v} is *not* unitary. We will return to this point later in Sect. 2.4.2. We have $\operatorname{div} \mathbf{v} = 1$ so that combining (6) and (7) yields

$$V = \int_{\Gamma} \langle \mathbf{v}, \mathbf{n} \rangle d\Gamma. \quad (8)$$

Let us look at the discrete case: Here, a level set Γ_h of (3) is extracted by the marching cubes in the form of a triangular mesh [14]. Such a piece-wise linear representation of the geometry provides a likewise locally-linear approximation of any function f whose values f_i are known at the vertices $\mathbf{x}_i \in \Gamma_h$, $i \in \mathbb{N}$. A Gauss-Legendre quadrature rule for f with linear precision defined over the triangle T is

$$\int_T f d\mathbf{x} \approx A(T) \sum_{i_k \in \mathcal{I}(T)} \frac{1}{3} f_{i_k},$$

where $A(T)$ equals the area of T , and $\mathcal{I}(T)$ enumerates its three corner vertices. Now substitute f_i by the flow $\langle \mathbf{v}_i, \mathbf{n}_i \rangle$. The vertex normals \mathbf{n}_i are usually taken

⁸But could be remedied by re-visiting the signed distance function φ of Γ from Sect. 2.3.

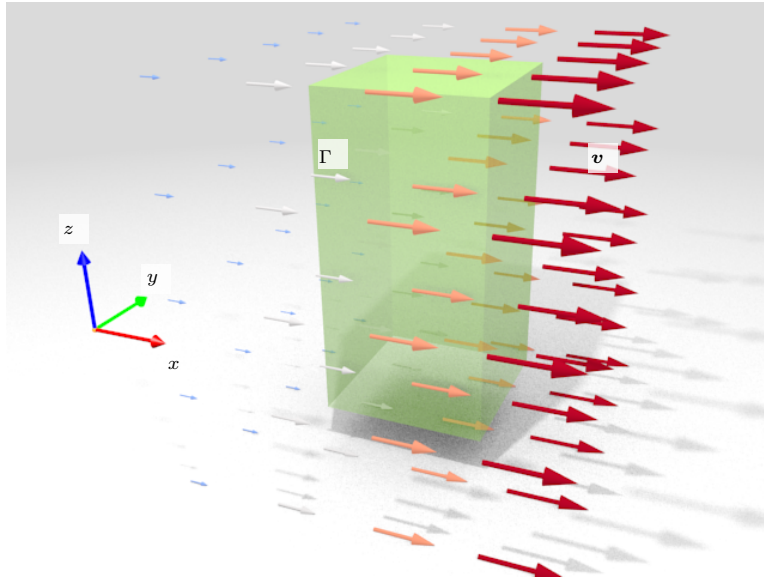


FIGURE 4. The volume of a compact object Ω equals the mean flow of the vector field $\mathbf{v} = (x, 0, 0)$ through its boundary surface Γ . If the support of the object is aligned with the xy - or xz -plane, it is not traversed by any stream line, and hence needs not be explicitly filled with mesh faces for a faithful volume estimate. The flow field is colored by increasing magnitude $\|\mathbf{v}\|$.

to be the normalized mean of the face normals in a one-ring neighborhood of \mathbf{x}_i . Altogether, we finally obtain the following approximation of Eq. (8)

$$V \approx \sum_{T \in \Gamma_h} \frac{A(T)}{3} \sum_{i_k \in \mathcal{I}(T)} \langle \mathbf{v}_{i_k}, \mathbf{n}_{i_k} \rangle.$$

2.4.2. Surfaces with boundary. As explained in Sect. 2.3, Poisson reconstruction accounts for most smaller holes in the aligned point clouds. The support of the object, i.e., its “bottom” or the area where it is in contact with the ground plane, however, remains usually unfilled. At the beginning of Sect. 2.4.1, we demanded that Γ be compact and closed because only then the volume of Ω is well-defined. We can lift this assumption in parts simply by a coordinate transformation: Remember that we chose \mathbf{v} to point in the direction of the x -axis of the world coordinate system. Consequently, the flow through any of the planes $y = \text{const}$ or $z = \text{const}$ vanishes. As shown in Fig. 4, all we have to do is align the support of the model with one of these planes. Without loss of generality, we choose $\{(x, y, z)^T \in \mathbb{R}^3 \mid z = 0\}$. In our scanning scenario, it is reasonable to assume that the specimens to be measured are spatially isolated enough that the depth images capture a significant portion of the ground plane surrounding the object, which can thus be detected fully automatically. To this end, we again invoke a RANSAC-type procedure, which samples triplets of points, calculates their common plane as a putative solution, and evaluates each such hypothetical plane by how many other points in the cloud it contains.

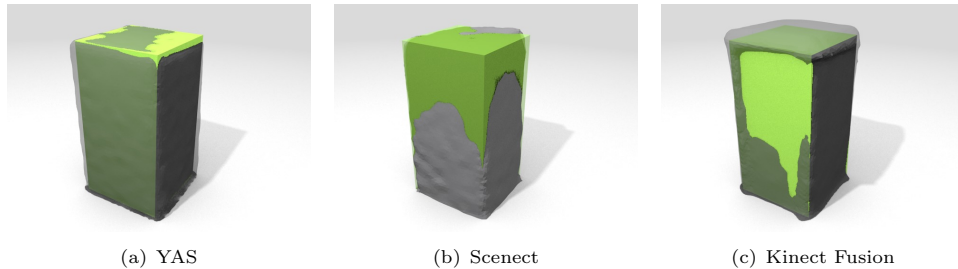


FIGURE 5. Reconstructions of item no. 1 in the cube data set. The ground truth is shown in green. It becomes apparent that Kinect Fusion Explorer systematically overestimates the object volume.

3. EXPERIMENTS

3.1. Implementation. We created a C++ implementation of most of the reconstruction pipeline, including raw data acquisition, coarse and fine registration as well as detection of the ground plane. Ease of use is of premier priority in view of the interdisciplinary nature of this project. Therefore, the number of dependencies was kept as small as possible: The OpenCV library supplies us with all functionality for feature matching (Sect. 2.2.1). Our implementation of the ICP method requires fast nearest-neighbor lookup which is based on the kd-tree data structure from the ANN library. We also created a graphical QT frontend which is showcased in the video included in the supplemental material. Poisson reconstruction (Sect. 2.3) is currently done in Meshlab but will be integrated into our code in upcoming releases.

We compare our method to the Kinect Fusion algorithm and Scenect, which is one of the few commercial software packages without hindering functionality restraints in the trial version. To warrant a fair comparison, all participating systems should be operated with the same sensor and the same intrinsic calibration. This proved to be somewhat difficult: Both Scenect and YAS access devices through the OpenNI framework driver supporting all of Primesense’s products and the Xtion by Asus among others. The Point Cloud Library provides an open-source version of the Kinect Fusion algorithm which could potentially function with the Carmine 1.09 as well, but our experiences with it were little encouraging. For this reason, we had to resort to Microsoft’s own implementation Kinect Fusion Explorer, which works exclusively with proprietary hardware, the Kinect.

3.2. Results. The set of 20 specimens can be divided into two equally-sized groups: The first group contains objects of simple geometry like cubes and cylinders, whose basic dimensions can be measured manually with a tape measure or ruler, see Fig. 2. Free-forms of sizes ranging from just a few centimeters (fifth row of Fig. 7) to the size of a human upper body (first row of Fig. 7) make up the second group. The “ground truth” volumes for the first group are listed in the first column of Tab. 1. Needless to say these are afflicted by their own uncertainty, given that they were determined through measurements with a ruler. Our reconstructions of the cube data set are depicted in Fig. 2. We obtain the best average relative volume error $(V_{\text{YAS}} - V_{\text{Manual}})/V_{\text{Manual}}$ of -0.34% . The performance of Scenect is comparable, which is somewhat surprising in view of Fig. 7(c). The meshes created by the Kinect

No.	V_{manual} [$10^{-3}m^3$]	V_{YAS} [$10^{-3}m^3$]	E_{YAS} [%]	V_{Scenect} [$10^{-3}m^3$]	E_{Scenect} [%]	V_{KinFu} [$10^{-3}m^3$]	E_{KinFu} [%]
1	1.26	1.26	0.11	1.26	0.6	1.47	16.6
2	2.82	2.72	-3.38	2.68	-5.1	3.61	28.4
3	1.29	1.21	-6.07	1.21	-6.17	1.63	27.0
4	3.07	3.07	-0.02	2.83	-7.57	4.23	38.0
5	2.18	2.06	-5.44	2.1	-3.79	2.85	31.1
6	6.26	6.44	2.86	5.76	-7.07	7.69	22.9
7	1.66	1.79	8.22	1.75	5.65	1.67	0.55
8	1.78	1.89	6.0	1.67	-6.28	2.15	20.3
9	2.75	2.75	0.05	2.68	2.41	3.84	39.8
10	24.6	24.8	1.05	18.6	-24.4	28.8	17.1
μ			-0.34		-5.74		24.2
σ			4.59		7.76		11.5

TABLE 1. Volume estimates and their relative error w.r.t. the value obtained by manual measurement.

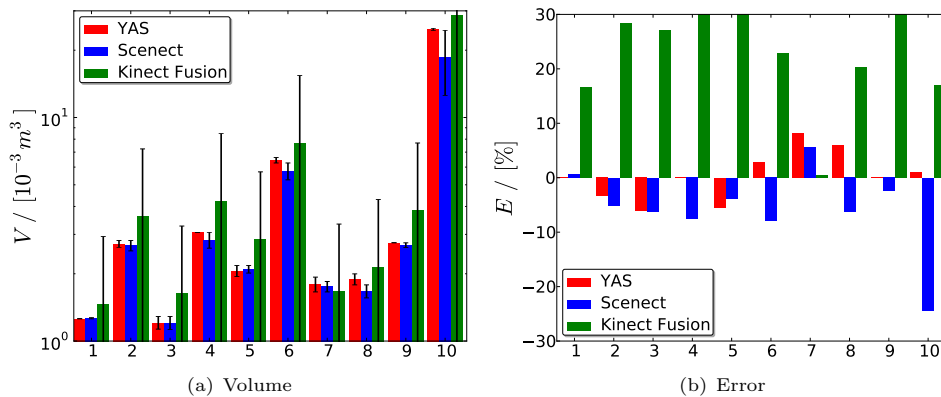


FIGURE 6. Visualization of Tab. 1. The error bars in (a) indicate the absolute deviation from the ground truth volume. From (b), it seems like Kinect Fusion is biased towards excessive values.

Fusion explorer are of inferior topological quality: they contain a high number of non-manifold simplices. This, however, does not seem to affect the volume estimates negatively. Also it can be said that the sensitivity of volumes w.r.t. the ground plane parameters is relatively low.

As can be seen from the last two columns of Tab. 1, the Kinect Fusion explorer systematically overestimates the ground truth volume by a significant margin. We conjecture that the issue is rooted in calibration. In fact, an important lesson learned during our experimental studies was that a good calibration can make a difference in error of an order of magnitude. Indeed, Scenect provides a calibration program, but Kinect Fusion Explorer does not. A visualization of Tab. 1 is plotted in Fig. 6.

Results for the second group of objects are shown in Fig. 7. Scenect performs worst among the three compared methods. It must be said, though, that Scenect

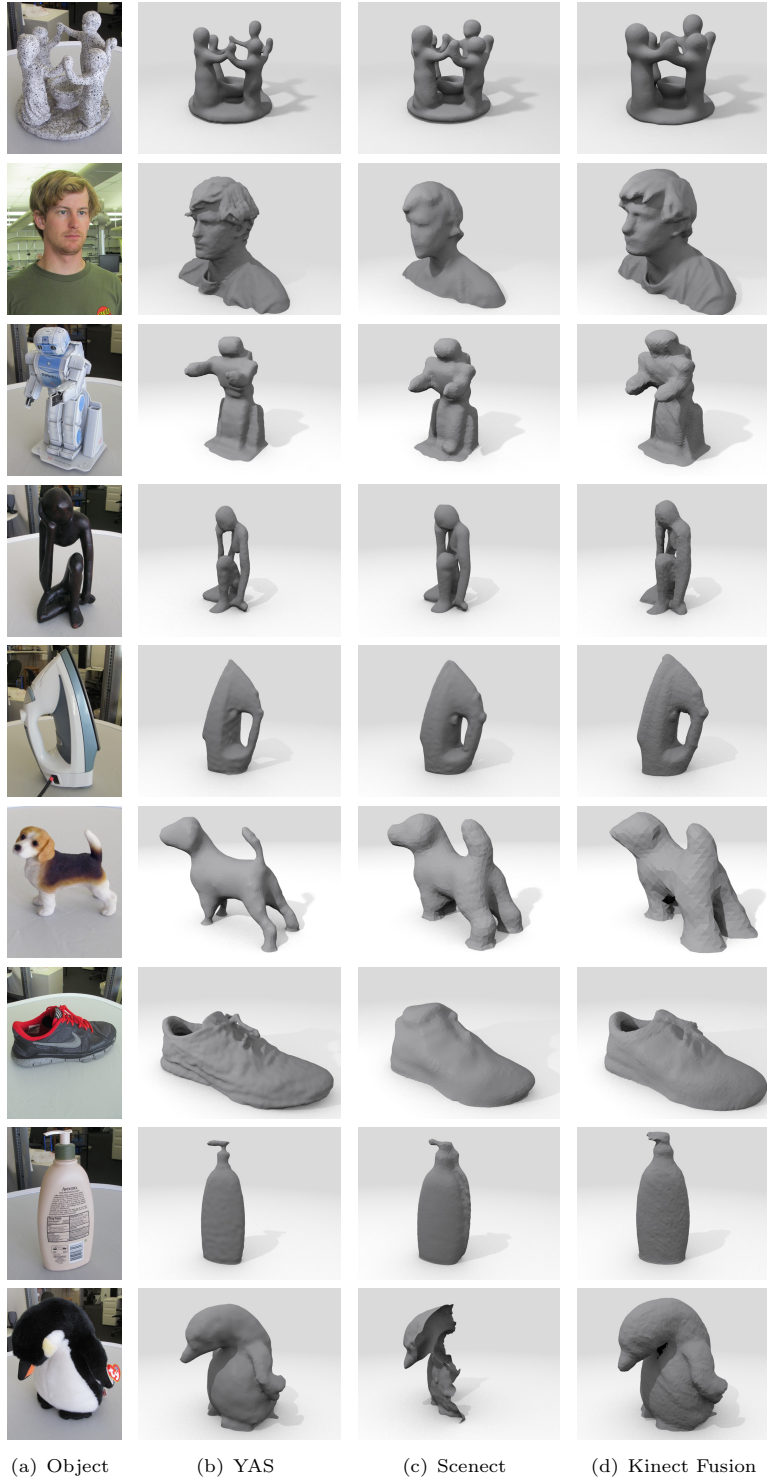


FIGURE 7. Selection of reconstructed free-forms. Let us emphasize that – to comply with the double-blind review policy – the bust in the first row is *not* taken from any one of the authors.

does not offer mesh reconstruction feature, and the point clouds it exports are not oriented. Normals can be computed by singular value decomposition considering the nearest neighbors of a point, which is probably less accurate than the finite-differences approximations of (5). The lack of loop-closure whose effect can be better identified in the point cloud in Fig. 3(b) carries over to the triangular mesh in row of Fig. 7(c). Premature termination of the tracker is responsible for the poor reconstruction of the penguin in the last row of Fig. 7(c).

Although it behaved unreliably during volume estimation, Kinect Fusion produced high quality models in real-time, which justifies the tremendous success it had since its inception. Still, the scale bias shows, and a lot of details appear to be missing – details which are present in our YAS reconstructions despite the fact that the Poisson algorithm is known to possess the characteristics of a low-pass filter (see the discussion in Sect. 2.3). The precision setting in Kinect Fusion is quite rigid since the dimensions of the Cartesian grid have to be fixed *before* the reconstruction process even starts. We found that the maximal resolution of $512 \times 512 \times 512$ voxels rendered the tracking unstable and/or led to incomplete meshes.

3.3. Discussion. We establish correspondence based on photometry, hence our approach fails in the absence of sufficiently exciting texture. This however does not necessarily need to be on the target object itself but can be found in the background, which may be “enriched” since after all, in our application, it is not of interest. The majority of competing approaches, including the two we compare against, can cope with the issue. The main reason is the trade-off between a sparse sampling of viewpoints and aforementioned need for sufficient texture: All previous method implicitly leverage on geometry in the motion estimation stage, which is only made possible by the small baseline between adjacent frames, i.e., through tracking, the disadvantages of which have been discussed in Sect. 2.2. In fact, ICP measures the similarity between two points by their distance, hence endows each of them with a geometry descriptor, although a primitive one, too primitive to support wide-baseline matching. More suitable descriptors are available, but we are not yet considering them here, for already the process of salient point detection let alone the computation of informative geometry descriptors is extremely challenging on noisy, occlusion-ridden, and incomplete depth data such as from RGBD sensors.

The system has no real-time capabilities. We believe this is not necessary for our target application of small-scale object scanning (unlike e.g. navigation, map-building, or odometry). Considering the low resolution of RGB images, the matching process is a matter of seconds. This bottleneck could be removed by an approximate nearest-neighbors search. Our impression is that our system requires the same or even less overall acquisition time compared to e.g. Kinect Fusion, which requires slow and steady motion around the object (and sometimes even a complete reset).

4. CONCLUSION

We described a system for integrating a set of RGBD images of small-scale objects into a geometrically faithful 3-d reconstruction. The system is intended to support researchers in the field of cognitive neuroscience who will use it for acquiring ground truth data for their own experimental studies. To assess the suitability of our 3-d models and those obtained by comparable algorithms, we performed an in-depth analysis of theoretical and empirical kind. There are two main conclusions

we would like to draw here: First, the quality of a set of calibration parameters or metric reconstruction can be deceptive. Only quantitative analysis enables veridical information. Second, while the uncertainties of hand-held structured-light scanners may be tremendous in the eyes of the optical metrologist, they help improving studies in the cognitive neurosciences, where manual measurement is still common practice.

REFERENCES

- [1] P. Besl and H. McKay. A method for registration of 3-D shapes. *IEEE T. Pattern Anal.*, 14(2):239–256, 1992.
- [2] J. Cantarella, D. DeTurck, and H. Gluck. Vector calculus and the topology of domains in 3-space. *Am. Math. Mon.*, 109(5):409–442, 2002.
- [3] A. Charpentier. Analyse expèrimentale: De quelques elements de la sensation de poids. *Archives of Physiology and Normal Pathology*, 3:122–135, 1891.
- [4] M. Ernst and M. Banks. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415:429–433, 2002.
- [5] M. Fischler and R. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, 1981.
- [6] J. Flanagan and M. Beltzner. Independence of perceptual and sensorimotor predictions in the size-weight illusion. *Nature Neurosci.*, 3:737–741, 2000.
- [7] B. Frayman and W. Dawson. The effect of object shape and mode of presentation on judgments of apparent volume. *Percept. Psychophys.*, 29:56–62, 1981.
- [8] M. Goodale. Transforming vision into action. *Vision Res.*, 51(13):1567–1587, 2011.
- [9] M. Grandy and D. Westwood. Opposite perceptual and sensorimotor responses to a size-weight illusion. *J. Neurophysiol.*, 95:3887–3892, 2006.
- [10] P. Henry, M. Krainin, E. Herbst, X. Ren, and D. Fox. RGB-D mapping: Using Kinect-style depth cameras for dense 3D modeling of indoor environments. *Int. J. Robot. Res.*, 31(5):647–663, 2012.
- [11] D. Herrera, J. Kannala, and J. Heikkilä. Joint depth and color camera calibration with distortion correction. *IEEE T. Pattern Anal.*, 34(10):2058–64, 2012.
- [12] M. Kazhdan, M. Bolitho, and H. Hoppe. Poisson Surface Reconstruction. *Eurographics SGP*, 1:61–70, 2006.
- [13] M. Kazhdan and H. Hoppe. Screened poisson surface reconstruction. *ACM Trans. Graph.*, 32(1):1–13, 2013.
- [14] W. Lorensen and H. Cline. Marching cubes: A high resolution 3D surface construction algorithm. *ACM SIGGRAPH '87*, 21(4):163–169, 1987.
- [15] D. Lowe. Object recognition from local scale-invariant features. *Proc. ICCV IEEE*, 1:1150–1157, 1999.
- [16] B. Mirtich. Fast and Accurate Computation of Polyhedral Mass Properties. *Journal of Graphics Tools*, 1(2):31–50, 1996.
- [17] R. Newcombe, D. Molyneaux, D. Kim, P. Koli, A. Davison, J. Shotton, S. Hodges, and A. Fitzgibbon. KinectFusion: Real-Time Dense Surface Mapping and Tracking. *IEEE Proc. ISMAR*, 1:127–136, 2011.
- [18] I. Rock and J. Victor. Vision and Touch: An Experimentally Created Conflict between the Two Senses. *Science*, 143(3606):594–596, 1964.
- [19] G. K. L. Tam, Z.-Q. Cheng, Y.-K. Lai, F. C. Langbein, Y. Liu, D. Marshall, R. R. Martin, X.-F. Sun, and P. L. Rosin. Registration of 3D point clouds and meshes: a survey from rigid to nonrigid. *IEEE T. Vis. Comput. Gr.*, 19(7):1199–217, 2013.
- [20] J. Williams and M. Bennamoun. Simultaneous registration of multiple corresponding point sets. *Comput. Vis. Image Und.*, 81(1):117–142, 2001.