

The Theory and Practice of Coplanar Shadowgram Imaging for Acquiring Visual Hulls of Intricate Objects

Shuntaro Yamazaki · Srinivasa G. Narasimhan ·
Simon Baker · Takeo Kanade

Received: 14 November 2007 / Accepted: 4 August 2008 / Published online: 20 September 2008
© Springer Science+Business Media, LLC 2008

Abstract Acquiring 3D models of intricate objects (like tree branches, bicycles and insects) is a challenging task due to severe self-occlusions, repeated thin structures, and surface discontinuities. In theory, a shape-from-silhouettes (SFS) approach can overcome these difficulties and reconstruct visual hulls that are close to the actual shapes, regardless of the complexity of the object. In practice, however, SFS is highly sensitive to errors in silhouette contours and the calibration of the imaging system, and has therefore not been used for obtaining accurate shapes with a large number of views. In this work, we present a practical approach to SFS using a novel technique called *coplanar shadowgram imaging* that allows us to use dozens to even hundreds of

views for visual hull reconstruction. A point light source is moved around an object and the shadows (silhouettes) cast onto a single background plane are imaged. We characterize this imaging system in terms of image projection, reconstruction ambiguity, epipolar geometry, and shape and source recovery. The coplanarity of the shadowgrams yields unique geometric properties that are not possible in traditional multi-view camera-based imaging systems. These properties allow us to derive a robust and automatic algorithm to recover the visual hull of an object and the 3D positions of the light source simultaneously, regardless of the complexity of the object. We demonstrate the acquisition of several intricate shapes with severe occlusions and thin structures, using 50 to 120 views.

This is an extension and consolidation of our previous work on coplanar shadowgram imaging system (Yamazaki et al. 2007) presented at IEEE International Conference on Computer Vision 2007.

Electronic supplementary material The online version of this article (<http://dx.doi.org/10.1007/s11263-008-0170-4>) contains supplementary material, which is available to authorized users.

S. Yamazaki (✉)
National Institute of Advanced Industrial Science and
Technology, Tokyo, Japan
e-mail: shun-yamazaki@aist.go.jp

S.G. Narasimhan · T. Kanade
Carnegie Mellon University, Pittsburgh, PA, USA

S.G. Narasimhan
e-mail: srinivas@cs.cmu.edu

T. Kanade
e-mail: tk@cs.cmu.edu

S. Baker
Microsoft Research, Redmond, WA, USA
e-mail: sbaker@microsoft.com

Keywords Multi-view geometry · Shape reconstruction · Shape from silhouette · Imaging system · Calibration · Intricate shape · Shadowgram · Coplanar shadowgram

1 Introduction

Acquiring 3D shapes of objects that have numerous occlusions, discontinuities and repeated thin structures is challenging for vision algorithms. For instance, the wreath object shown in Fig. 1(a) contains over 300 branch-lets each 1–3 mm in diameter and 20–25 mm in length. Covering the entire surface area of such objects requires a large number (dozens or even a hundred) of views. Thus, finding correspondences between views as parts of the object get occluded and “dis-occluded” becomes virtually impossible, often resulting in erroneous and incomplete 3D models.

The issues of correspondence and occlusion in the object can be avoided if we only use the silhouettes of an

Fig. 1 Obtaining 3D models of intricate shapes such as in (a) is hard due to severe occlusions and correspondence ambiguities. (b) By moving a point source in front of the object, we capture a large number of shadows cast on a single fixed planar screen (122 views for this object). Applying our techniques to such coplanar shadowgrams enables the accurate recovery of intricate shapes

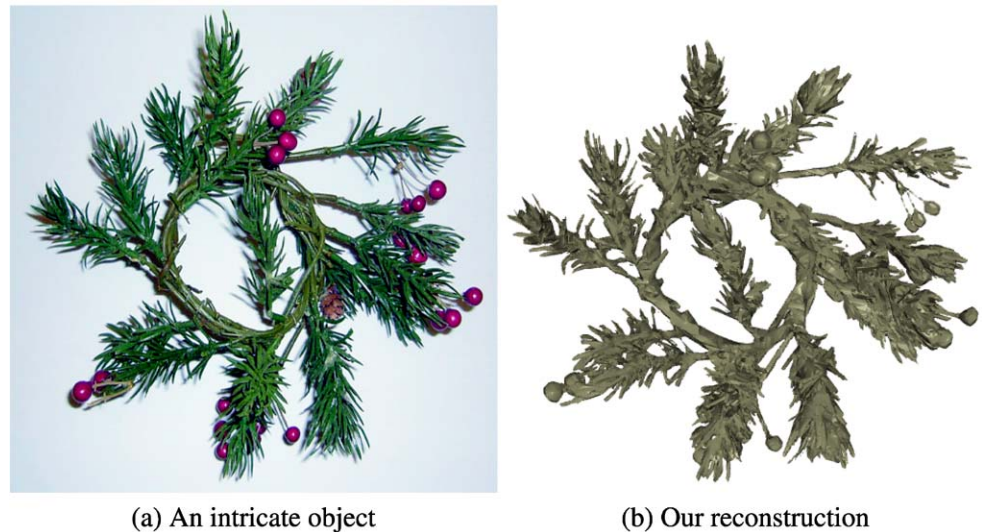
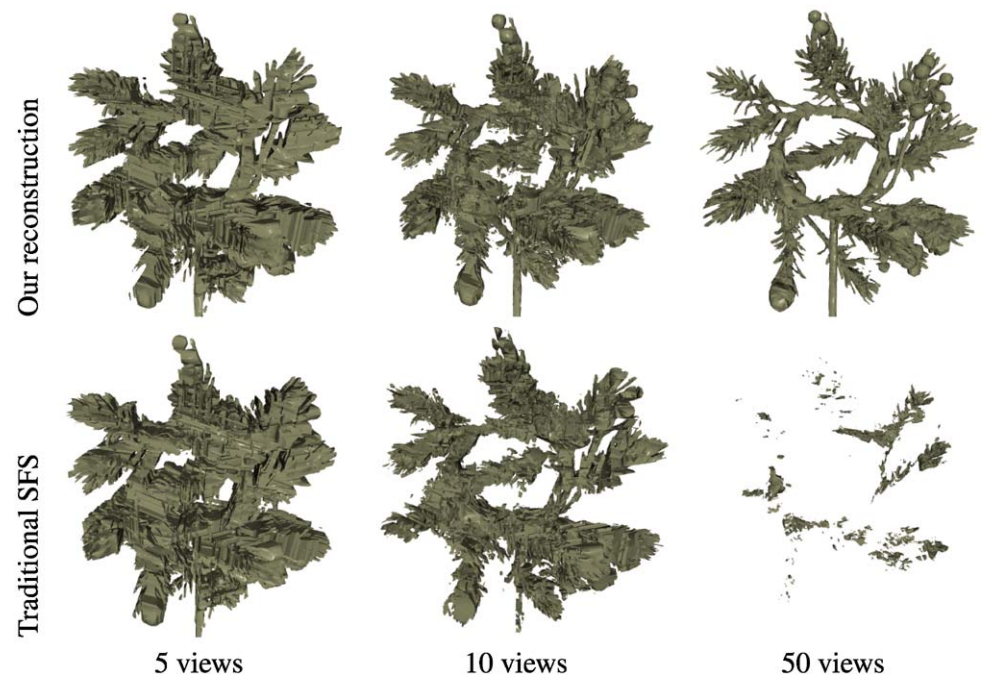


Fig. 2 Sensitivity of SFS reconstruction: (Top) The visual hulls reconstructed using the light source positions estimated by our method. As the number of silhouettes increases, the visual hull gets closer to the actual shape. (Bottom) The reconstructions obtained from slightly erroneous source positions. As the number of views increases, the error worsens significantly



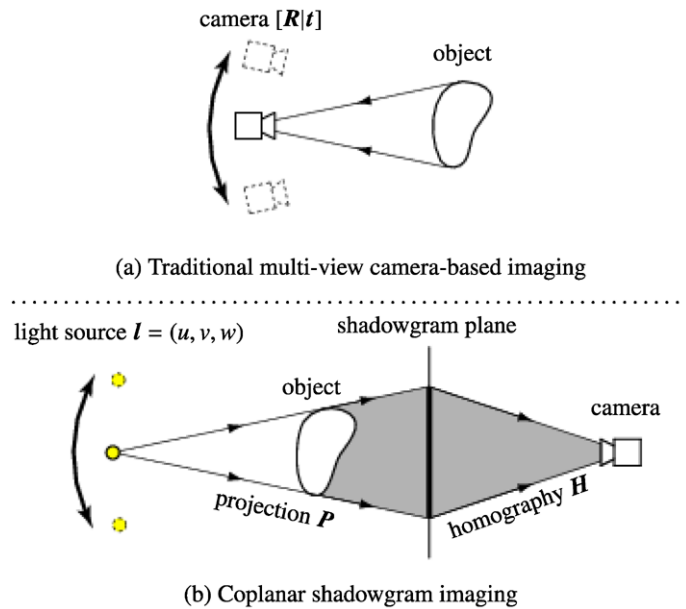
object obtained from different views and reconstruct its *visual hull* (Baumgart 1974). The top row of Fig. 2 illustrates the visual hulls estimated using our technique from different numbers of silhouettes. While the visual hull computed using a few (5 or 10) silhouettes is too coarse, the reconstruction from a large number of views (50) is an excellent model of the original shape. Thus, a Shape-From-Silhouettes (SFS) approach may be used to recover accurate 3D models of intricate objects.

In practice, however, SFS algorithms are highly sensitive to errors in silhouette contours and the geometric parameters of the imaging system (camera calibration) (Sinha et al. 2004). This sensitivity worsens as the number of views in-

creases, resulting in poor quality models. The bottom row in Fig. 2 shows the visual hulls of the wreath object obtained from slightly erroneous source positions. This drawback must be addressed in order to acquire intricate shapes reliably.

The traditional SFS system assumes that a camera observes the object, and the silhouette is extracted from the obtained image by image processing, such as, background subtraction (Cheung et al. 2005; Smith and Blinn 1996), image segmentation (Boykov and Funka-Lea 2006; Campbell et al. 2007) or manual segmentation (Furukawa et al. 2006). Multiple viewpoints are captured by moving either the camera or the object (see Fig. 3(a)). For each view, the relative

Fig. 3 (a) The object of interest is observed directly by a camera. The silhouette of the object is extracted from the captured image. Multiple views are obtained by moving the camera or the object. (b) A point source illuminates the object and its shadow cast on a planar rear-projection screen represents the silhouette of the object. Coplanar shadowgrams from multiple viewpoints are obtained by moving the light source. Note that the relative transformation between the object and the screen remains fixed across different views



pose between the object and the camera is described by six parameters (3D translation and 3D rotation). Savarese et al. (2005) proposed a system that avoids the difficulty in silhouette extraction using cast shadows. When an object is illuminated by a single point light source, the shadow cast onto a background plane is sharp and can be directly used as its silhouette. Silhouettes from multiple views are obtained by rotating the object and capturing the shadow images (also known as *shadowgrams*¹). Balan et al. (2007) proposed a hybrid method that reconstructs the visual hull of an object using the shadows cast on the floor as well as the silhouettes captured by multiple cameras. In terms of multi-view geometry, these methods are equivalent to traditional SFS, requiring six parameters per view.

This paper proposes a novel approach to SFS called *coplanar shadowgram imaging*. We use a setup similar in spirit to that proposed by Savarese et al. (2005). The key difference here is that the point source is moved, while the object, the camera and the background screen all remain stationary. The principal focus of this work is on acquiring visual hulls of intricate and opaque objects from a large number of coplanar shadowgrams. Our main contributions are described below.

Multi-View Geometry of Coplanar Shadowgram Imaging: We propose a coplanar shadowgram imaging system, and clarify how it is different from traditional imaging system where the camera directly observes the silhouette of an object. Figure 3 compares the geometry of the two systems.

¹Shadowgrams have also been widely used for visualizing the 3D structure of transparent objects such as glasses or fluids (Settles 2001; Hooke 1667).

The key observation is that the relative transformation between the object and screen remains fixed across different views in our coplanar shadowgram imaging system. The image projection model is hence described by only three parameters per view (3D translation of the source) instead of six in traditional systems. The geometry is similar in spirit to parallax geometry (Sawhney 1994; Cross et al. 1999) where the homography between image planes is known to be an identity, which allows the derivation of unique geometric properties that are not possible in the traditional multi-view camera-based imaging system. For instance, we show that epipolar geometry can be uniquely estimated from only the shadowgrams, without requiring any correspondences, and independent of the object's shape.

Recovery of Shape and Source Positions: When the shape of the object is unknown, the locations of all the point sources (and therefore the object's shape) can be recovered from coplanar shadowgrams, only up to a four parameter perspective transformation. We show how this transformation relates to the Generalized Perspective Bas-Relief (GPBR) ambiguity (Kriegman and Belhumeur 2001) that is derived for a single viewpoint system. We break this ambiguity by simultaneously capturing the shadowgrams of two or more spheres and using them as soft constraints.

Robust Reconstruction of Visual Hull: Even a small amount of blurring in the shadow contours may result in erroneous estimates of source positions that in turn can lead to erroneous visual hulls due to the non-linear nature of the reconstruction algorithm. We propose an optimization of the light source positions that can robustly reconstruct the visual hulls of intricate shapes. First, the large error in light

source positions is corrected by enforcing the reconstructed epipolar geometry. We then minimize the mismatch between the acquired shadowgrams and those obtained by reprojecting the estimated visual hull. Undesirable local convergence in the non-linear optimization is alleviated using the convex polygons of the silhouette contours. In practice, the optimization on the convex silhouettes also leads to faster convergence.

For the analogous camera-based imaging, a number of algorithms have been proposed to make SFS robust to errors in camera position and orientation. These techniques optimize camera parameters by exploiting either epipolar tangency or silhouette consistency.

When we exploit the epipolar tangency and reconstruct the projective geometry of imaging system, the localization of epipole and the correspondence between tangent lines have to be solved simultaneously. To solve this combinatorial optimization, existing methods assume either good initializations (Cipolla et al. 1995; Sinha et al. 2004) or an approximated camera model such as linear projection (Furukawa et al. 2006). These methods also assume that the shape of the object is *reasonably complex*; Silhouettes of simple objects such as spheres do not have enough features, and intricate objects like branches have too many, making it hard to find correspondences automatically. Wong and Cipolla propose a technique that can recover the epipolar geometry regardless of the complexity of the silhouette, using the outer epipolar tangency which is similar in spirit to our approach (Wong and Cipolla 2004). The camera motion, however, is limited to circular motion since they use the multi-view camera-based imaging system.

Another approach to the camera reconstruction from silhouettes is estimating the camera poses that can reconstruct the visual hull consistent to with the acquired silhouettes. Yezzi and Soatto (2003) solved this highly non-linear optimization problem assuming a good initialization and simple silhouettes. Hernández et al. proposed an efficient optimization by maximizing the consistency only along the contour lines of the silhouettes. However, the methods based on the silhouette consistency also have difficulty in the initialization and the limitation on the complexity of the shape. As a result, these techniques have succeeded in only acquiring the 3D shapes of relatively simple objects like people and statues using a small number of views and/or assuming limited camera motion.

In contrast, our algorithm is effective for a large number of views (dozens to a hundred) from a wide range of viewpoints, does not require any feature correspondences and does not place any restriction on the shapes of the objects. The minimization of silhouette mismatch is also easier requiring optimization of source translation (3 DOF per view), instead of the harder (and sometimes ambiguous (Hartley and Zisserman 2004)) joint estimation of camera rotation

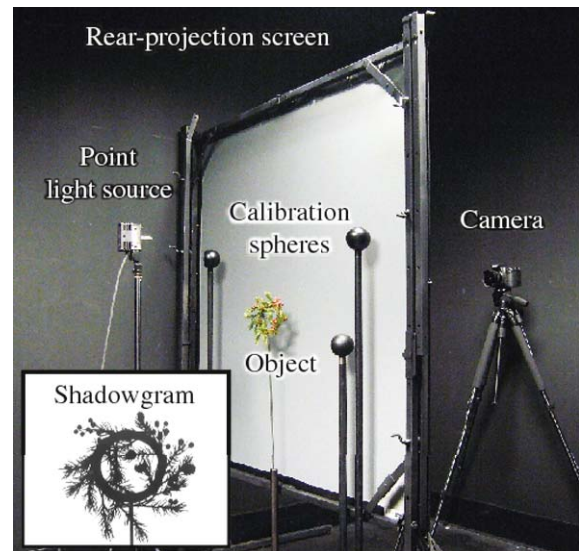


Fig. 4 The setup used to capture coplanar shadowgrams includes a single point light source, a rear-projection screen, and a digital camera. The object is placed close to the screen to cover a large field of view. Two or more spheres are used to estimate the initial light source positions. (Inset) An example shadowgram obtained using the setup

and translation (6 DOF per view) in the traditional system. As a result, we achieve good quality reconstructions of real objects such as a wreath, a wiry ball, a palm tree, and an octopus that show numerous occlusions, discontinuities and thin structures. In addition, we have also evaluated our techniques quantitatively using simulations with objects such as a coral, branches, a bicycle and an insect whose ground truth shapes are known beforehand.

Despite significant progress in optical scanning hardware (Curless and Levoy 1996; Levoy et al. 2000) and multi-view geometry (Hartley and Zisserman 2004; Seitz et al. 2006), reconstruction of intricate shapes remains an open problem. We believe this work is an initial step to solving this problem. Possible extensions of our work include multiple screens covering $360^\circ \times 360^\circ$ views of the objects and/or multiple light sources for dynamic objects, and combine our techniques with stereo and photometric stereo, to obtain reconstructions that are smoother than visual hulls, including concavities.

2 Coplanar Shadowgrams

We define *shadowgrams* as the shadow pattern cast on a background plane by an object that occludes a point source. If the object is opaque, the sharp shadowgram accurately represents the silhouette of the object. Henceforth, we shall use shadowgrams and silhouettes interchangeably. We also refer to the images that captures shadowgrams as shadowgram images. Coplanar shadowgram imaging is the process

of acquiring several shadowgrams on a *single plane* by moving the light source while keeping the shadow plane fixed.

Our acquisition setup shown in Fig. 4 includes a 6M-pixel Canon EOS-20D digital camera, a 250 watt 4 mm incandescent bulb, and a 4 ft \times 4ft translucent rear-projection screen. This setup is similar to the one in Savarese et al. (2005) where the object is rotated and the light source is fixed. We prove in Sect. 4 that the fixed relative position between the image plane and the object is crucial to deriving a strong constraint on the shadowgrams. Another advantage of moving only the light source is that we can easily acquire the shadowgrams of non-rigid objects. In this section, we describe how the visual hull of the object is obtained using our setup.

2.1 Shadowgram Projection

Throughout this paper, we represent the locations of light sources, the surface of 3D objects, and the shadowgrams on a background plane in either a two or three dimensional real projective space (\mathbb{RP}^2 or \mathbb{RP}^3) unless defined explicitly. Sets are written in calligraph face as \mathcal{S} , operators are in sans-serif font as P , vectors and matrices are in bold face as \mathbf{v} and \mathbf{M} , and scalars are in normal fonts as x . All the vectors without the transpose are column vectors.

Figure 3(b) illustrates the viewing and illumination geometry of coplanar shadowgram imaging. Without loss of generality, we assume the shadowgram plane Π is located at $z = 0$. When an 3D object $\mathcal{O} \subset \mathbb{RP}^3$ is illuminated by a point light source at $\mathbf{l} \in \mathbb{RP}^3$, the shadowgram $\mathcal{S}^{\mathcal{O}} \subset \mathbb{RP}^2$ of the object on the shadowgram plane is obtained by a perspective projection:

$$\mathcal{S}^{\mathcal{O}} = P(\mathbf{l}) \cdot \mathcal{O}. \quad (1)$$

The projection $P(\mathbf{l}) : \mathbb{RP}^3 \rightarrow \mathbb{RP}^2$ can be represented as a 3×4 matrix $\mathbf{P}(\mathbf{l})$ (see Appendix A for the derivation):

$$\mathbf{P}(\mathbf{l}) = \begin{pmatrix} -w & 0 & u & 0 \\ 0 & -w & v & 0 \\ 0 & 0 & 1 & -w \end{pmatrix}. \quad (2)$$

The image \mathcal{I} of the shadowgram acquired by a camera is related to the shadowgram $\mathcal{S}^{\mathcal{O}}$ on the plane Π by a 2D homography:

$$\mathcal{I} = \mathbf{H} \cdot \mathcal{S}^{\mathcal{O}}. \quad (3)$$

This homography $\mathbf{H} : \mathbb{RP}^2 \rightarrow \mathbb{RP}^2$ is independent of the light source position and can be estimated separately using a standard computer vision algorithm. In the following, we assume that the shadowgram $\mathcal{S}^{\mathcal{O}}$ has been estimated using

$$\mathcal{S}^{\mathcal{O}} = \mathbf{H}^{-1} \cdot \mathcal{I}. \quad (4)$$

2.2 Visual Hull Reconstruction

Now let a set of shadowgrams $\mathcal{S}_i^{\mathcal{O}}$ of an object \mathcal{O} ($i = 1, \dots, N$) be acquired by moving the source to different locations \mathbf{l}_i . The visual cone $\mathcal{C}_i \subset \mathbb{RP}^3$ associated with the shadowgram $\mathcal{S}_i^{\mathcal{O}}$ is defined as

$$\mathcal{C}_i \stackrel{\text{def}}{=} \left\{ \mathbf{p} \in \mathbb{RP}^3 \mid P(\mathbf{l}_i) \cdot \mathbf{p} \in \mathcal{S}_i^{\mathcal{O}} \right\}. \quad (5)$$

Then, the visual hull $\mathcal{V} \subset \mathbb{RP}^3$ of the object is obtained as

$$\mathcal{V} \stackrel{\text{def}}{=} \bigcap_{i=1}^N \mathcal{C}_i. \quad (6)$$

Given the 3D locations \mathbf{l}_i of the light sources, the visual hull of the object can be estimated using (2) and (6). Due to the nature of the reconstruction algorithm, the reprojection $\mathcal{S}^{\mathcal{V}}$ of the visual hull \mathcal{V} to the shadowgram plane reproduces the shadowgram identical to the input.

$$\mathcal{S}_i^{\mathcal{V}} \stackrel{\text{def}}{=} P(\mathbf{l}_i) \cdot \mathcal{V} \quad (7)$$

$$= \mathcal{S}_i^{\mathcal{O}}. \quad (8)$$

Table 1 summarizes and contrasts the geometric parameters that appear in the traditional multi-view camera-based and coplanar shadowgram imaging systems. In multi-view camera-based imaging system (Fig. 3(a)), we must specify both the intrinsic camera parameters that define the relationship between pixels and ray directions, and the extrinsic camera parameters (rotation and translation) that define the ray directions in a world coordinate frame. As the camera moves, the intrinsic parameters remain the same, while the six extrinsic parameters vary. On the other hand, in coplanar shadowgram imaging (Fig. 3(b)), the point source is defined only by its location, requiring three parameters per view as compared to the six in the traditional system. As we shall show, this difference is crucial to the success of the technique presented in this paper.

3 Source Recovery

When the shape of the object is unknown, it is not possible to uniquely recover the 3D source positions only using the coplanar shadowgrams. In this section, we discuss the nature of the ambiguity in the 3D reconstruction from our coplanar shadowgrams and propose a simple method to break it.

3.1 Ambiguity in 3D Reconstruction from Coplanar Shadowgrams

Suppose two sets of light source locations, $\mathbf{l}_i = (u_i, v_i, w_i, 1)^T$ and $\mathbf{l}'_i = (u'_i, v'_i, w'_i, 1)^T$ for $i = 1, \dots, N$,

are estimated from a given set of coplanar shadowgrams. We can reconstruct respective visual hulls \mathcal{V} and \mathcal{V}' using (6). The reprojection of the visual hulls defined in (7) satisfy:

$$P(I_i) \cdot \mathcal{V} = P(I'_i) \cdot \mathcal{V}'. \quad (9)$$

Let $A: \mathbb{RP}^3 \rightarrow \mathbb{RP}^3$ be a transformation that deforms \mathcal{V} into \mathcal{V}' :

$$\mathcal{V}' = A \cdot \mathcal{V}. \quad (10)$$

Equations (9) and (10) lead to the identical equation:

$$P(I_i) = P(I'_i) \cdot A. \quad (11)$$

Here, A is a linear transformation because any 3D point along a ray from a point light source at I_i is transformed to the point on a ray from the source at I'_i . Comparing all elements in both sides of (11), we obtain the matrix A of the transformation (see Appendix B for the derivation):

$$A = \begin{pmatrix} 1 & 0 & a_1 & 0 \\ 0 & 1 & a_2 & 0 \\ 0 & 0 & a_3 & 0 \\ 0 & 0 & a_4 & 1 \end{pmatrix}. \quad (12)$$

The four parameters a_1, a_2, a_3 , and a_4 take arbitrary numbers under the condition that

$$a_3 \neq 0 \quad \text{and} \quad a_3 a_4 \geq 0. \quad (13)$$

Using (11), we can prove that the transformation between the source locations is also described by the same matrix:

$$I'_i = A I_i. \quad (14)$$

Equations (10) and (12) show that we can recover the shape of an object only from coplanar shadowgrams, up to a four-parameter family of perspective transformations A .

This result is consistent with the theory on the ambiguity in the geometric reconstruction from shadow information. Kriegman and Belhumeur (2001) prove that, when the shape of an object and the locations of light sources are all unknown, we can at best reconstruct the 3D structure from the shadow information, up to a four-parameter family of perspective transformation which they call the Generalized Perspective Bas-Relief (GPBR) transformation. In fact, our ambiguity transformation A can be viewed as the GPBR transformation if the origin of the world coordinate system is translated to the center of projection of a calibrated camera, and the shadowgram plane is regarded as a part of the object. In our imaging system, however, we use the coordinate system where the origin is on the shadowgram plane for the following reasons. Firstly, the estimation of the homography between the camera image plane and the shadowgram plane is much easier than the calibration of camera

projection. We can reconstruct the shape of the object without knowing any camera parameter. Secondly, the ambiguity in the global scale of the reconstruction can be resolved using the physical dimension of the shadowgram screen.

3.2 Determining Source Position Using Spheres

Significant effort has been devoted to understanding when and how the ambiguity in the shape reconstruction using shadows can be resolved. If the surface reflection on the object of interest is observed by a camera, the ambiguity is resolved by exploiting relative albedo values (Hayakawa 1994), interreflections (Chandraker et al. 2005), specularities (Georghiades 2003; Drbohlav and Sara 2002; Drbohlav and Chantler 2005), or Helmholtz reciprocity (Tan et al. 2007). However, we cannot observe the surface reflection in coplanar shadowgram imaging system. The robustness of the analyses of surface reflection is also open to question when the object has an intricate shape.

We now propose a technique of breaking the ambiguity using calibration objects. The location $I = (u, v, w, 1)^T$ of a point light source is directly estimated by capturing shadowgrams of two additional spheres that are placed adjacent to the object of interest.

The first and second coordinates u and v can be estimated by analyzing the shadowgrams of the spheres. Figure 5 illustrates the coplanar elliptical shadowgrams cast by the two spheres.² The ellipses are localized using a constrained least squares approach (Fitzgibbon et al. 1999). The intersection of the major axes $\overline{A_1 B_1}$ and $\overline{A_2 B_2}$ of the two ellipses is the foot of the perpendicular line from I to Π , denoted by $I_\perp = (u, v, 0, 1)^T$ in Fig. 5.

The third coordinate w is obtained as the intersection of hyperbolae in 3D space as shown below. Without loss of generality, consider the 3D coordinate system whose origin is at the center of the ellipse, and X and Y axes are respectively the major and minor axes of the ellipse. Then, the ellipse is represented in the following form.

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1 \quad (a > b). \quad (15)$$

In 3D space, there exists an inscribed sphere tangent to the conical surface and the plane, regardless of the position or the radius of the sphere. The cross section of the inscribed sphere by the plane that includes the apex of the cone and the major axis of the ellipse is shown in Fig. 6. The center of

²Each sphere is placed so that the minimum distance between a light source and the rear-projection screen is larger than the distance between the center of the sphere and the screen. Under this configuration, the cast shadow of the sphere is always an ellipse (Besant 1890).

the inscribed sphere is shown by R . The other symbols are corresponding to those in Fig. 5. The center of the ellipse C is the origin of the coordinate system.

The inscribed sphere is tangent to XY -plane at a focus of the ellipse R' , hence

$$\overline{CR'} = \sqrt{a^2 - b^2}. \quad (16)$$

Using the symmetry of triangles,

$$\overline{LA} - \overline{AR'} = \overline{LB} - \overline{BR'}. \quad (17)$$

Let the position of the apex be $\mathbf{l} = (t, 0, w, 1)^T$ in this coordinate system, then we can solve w with respect to t as:

$$w = \sqrt{\frac{b^2 t^2}{a^2 - b^2} - b^2}, \quad (18)$$

where a and b are the semimajor and semiminor axes of one of the ellipses, and t is the length between \mathbf{l}_\perp and the center of the ellipse.

More than two spheres may be used for a robust estimate of the source position. The above method is completely automatic and does not require the knowledge of the radii of the spheres, the exact locations at which they are placed in the scene, or point correspondences.

This technique for estimating the source position can be sensitive to errors in measured silhouettes. Due to the finite size of the light bulb, the shadowgram formed may be blurred, making it hard to localize the boundary of the silhouette. The extent of blurring depends on the relative distances of the screen and source from the object. To show the sensitivity of the technique, we performed simulations with spheres. We blurred the simulated silhouettes (effective resolution 480×360 pixels) with 5×5 and 10×10 averaging kernels, and estimated the 3D coordinates of the light source. Figure 7 presents u and w components of the source positions reconstructed using three spheres. Observe that the estimation becomes poor when the shadowgram is close to a right circle. In turn, the visual hull of a tree branch computed from the erroneous source positions is too erroneous even to perceive the 3D structure.

Due to the nature of visual hull reconstruction, a large error in one light source can blast the perfect reconstruction obtained by the other sources. Thus, better algorithms for averaging out the errors in individual light sources are crucial for obtaining accurate 3D models of intricate shapes.

4 Epipolar Geometry

Analogous to the scenario of binocular stereo, we define the epipolar geometry between a pair of shadowgrams that are generated by placing the point source in two locations

($\mathbf{l}_i = (u_i, v_i, w_i, 1)^T$ and $\mathbf{l}_j = (u_j, v_j, w_j, 1)^T$ in Fig. 8). Here, the locations of the point source are analogous to the centers-of-projection of the stereo cameras. The baseline connecting the two light sources \mathbf{l}_i and \mathbf{l}_j intersects the shadowgram plane Π at the epipole \mathbf{e}_{ij} . When the light sources are equidistant from the shadowgram plane Π , the epipole is at infinity.

Based on these definitions, we make two key observations that do not hold for binocular stereo: since the shadowgrams are coplanar, (a) they share the *same epipole* and (b) the points on the two shadowgrams corresponding to the same scene point lie on the *same epipolar line*. These observations are respectively written as

$$\mathbf{M}_{ij} \mathbf{e}_{ij} = 0, \quad (19)$$

$$\mathbf{m}_i^T \mathbf{F}_{ij} \mathbf{m}_j = 0. \quad (20)$$

In (19), \mathbf{M}_{ij} is a 2×3 matrix composed of two plane equations in the rows

$$\mathbf{M}_{ij} = \begin{pmatrix} -\Delta v & \Delta u \\ \Delta u \Delta w & -\Delta - v \Delta w \\ u_i v_j - u_j v_i \\ (u_i \Delta u + v_i \Delta v) \Delta w - w_i (\Delta u^2 + \Delta v^2) \end{pmatrix} \quad (21)$$

where $\Delta u = u_j - u_i$, $\Delta v = v_j - v_i$, and $\Delta w = w_j - w_i$. In (20),

$$\mathbf{F}_{ij} = [\mathbf{e}_{ij}]_\times \quad (22)$$

is the *fundamental matrix* that relates two corresponding points \mathbf{m}_i and \mathbf{m}_j between shadowgrams. $[\mathbf{e}_{ij}]_\times$ is the 3×3 skew symmetric matrix:

$$[\mathbf{e}_{ij}]_\times \mathbf{x} = \mathbf{e}_{ij} \times \mathbf{x} \quad (23)$$

for any 3D vector \mathbf{x} .

The camera geometry in coplanar shadowgram is similar in spirit to the parallax geometry (Sawhney 1994; Cross et al. 1999) where the image deformation is decomposed into a planar homography and a residual image parallax vector. In our system, however, the homography is exactly known to be an identity, which allows us to recover the epipolar geometry *only* from acquired images accurately regardless of the number of views or the complexity of the shadowgram contours.

4.1 Algorithm for Estimating Epipolar Geometry

We now show that the above observations enable us to uniquely estimate the epipolar geometry from only the shadowgram images. Suppose we have the plane in Fig. 8 that includes the baseline and is tangent to the surface of an object at a *frontier point* F . The intersection of this plane and

Fig. 5 Source position $\mathbf{l} = (u, v, w, 1)^T$ is recovered using the elliptical shadowgrams of two spheres. The radii and positions of the spheres are unknown. The major axes of the ellipses intersect the screen at $\mathbf{l}_\perp = (u, v, 0, 1)^T$. The w component is obtained using (18)

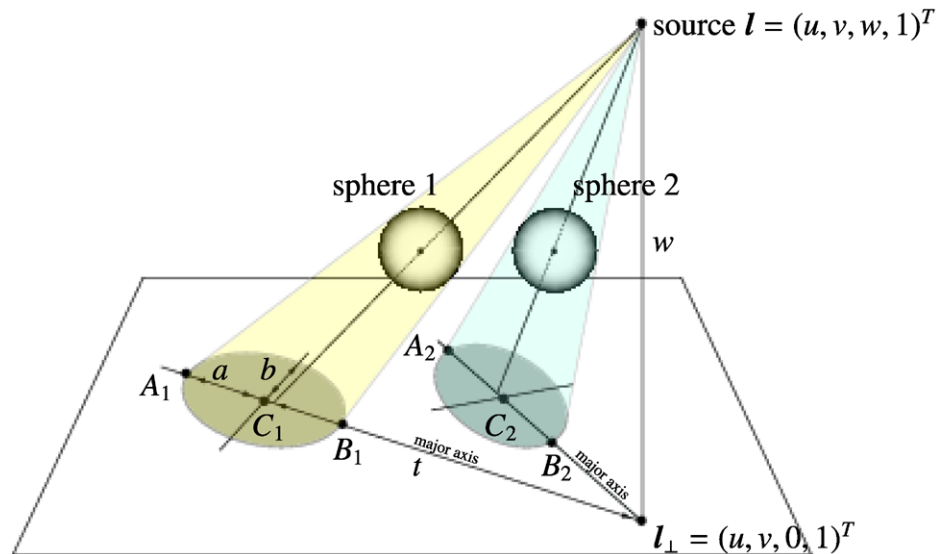


Table 1 Comparison between the geometric parameters of silhouette projection. For n views, the traditional multi-view system is described by $5 + 6n$ parameters. In comparison, the coplanar imaging system requires only $8 + 3n$ parameters

	View independent	View dependent
Projective cameras	1 (focal length)	3 (rotation)
	1 (aspect ratio)	3 (translation)
	1 (skew)	
	2 (image center)	
Coplanar shadowgrams	8 (homography H)	3 (translation \mathbf{l})

the shadowgram plane Π forms an epipolar line, which is also known as an *epipolar bitangent* (Cross et al. 1999), that can be estimated as one that is cotangent to the two shadowgrams (at T_i and T_j in Fig. 8). Two such epipolar lines can then be intersected to localize the epipole. But the reverse need not hold; every cotangent line need not be an epipolar line. So, how do we localize the epipole and estimate the epipolar lines without knowing any frontier points?

Figure 9(a) illustrates the simplest case of two convex shadowgrams partially overlapping each other. There are only two cotangent lines that touch the shadowgrams at the top and bottom region, resulting in a unique epipole e . When the convex shadowgrams do not overlap each other, four distinct cotangent lines are possible, generating six candidate epipoles, as shown by dots in Fig. 9(b). Only two of these four cotangent lines pass through the actual epipole, hence, the other two are false detections. Indeed, the false detections correspond to infeasible cases where the light source is located between the object and the screen, or behind the screen. We can detect actual epipolar lines by choosing the

cotangent lines where the epipole does not appear between the two points of shadowgram tangency.

When shadowgrams are non-convex, the number of cotangent lines can be arbitrarily large depending on the complexity of the shadowgram contours. Figure 9(c) illustrates the multiple candidates of cotangent lines at the point of tangency T . In this case, we compute the convex polygon surrounding the silhouette contour as shown in Fig. 9(d) and prove the following proposition (see Appendix C for the proof):

Proposition 1 *The silhouettes of the convex hull of an object are the convex hulls of the silhouettes.*

Using Proposition 1, the problem of estimating epipolar lines for concave silhouettes is reduced to the case of either (a) or (b). Thus, epipolar geometry can be reconstructed uniquely and automatically from only the shadowgrams. This capability of recovering epipolar geometry is independent of the shape of silhouette, and hence, the 3D shape of the object. Even when the object is a sphere, we can recover the epipolar geometry without any ambiguity. This geometric property is also observed in the parallel projection model where cameras are internally calibrated and moved without rotation (Åström et al. 1999). In the coplanar shadowgram imaging, however, the same geometric structure is accomplished without camera calibration or any apparatus to control the cameras. Table 2 summarizes the differences between traditional multi-view imaging by uncalibrated cameras and coplanar shadowgrams in terms of recovering epipolar geometry.

For the special case where the baseline intersects a convex object, one convex silhouette lies completely within the other and hence the epipole lies within the silhouettes. In

Fig. 6 The cross section of a right circular conical surface formed by the light rays emanating from a point light source at l and tangent to a calibration sphere

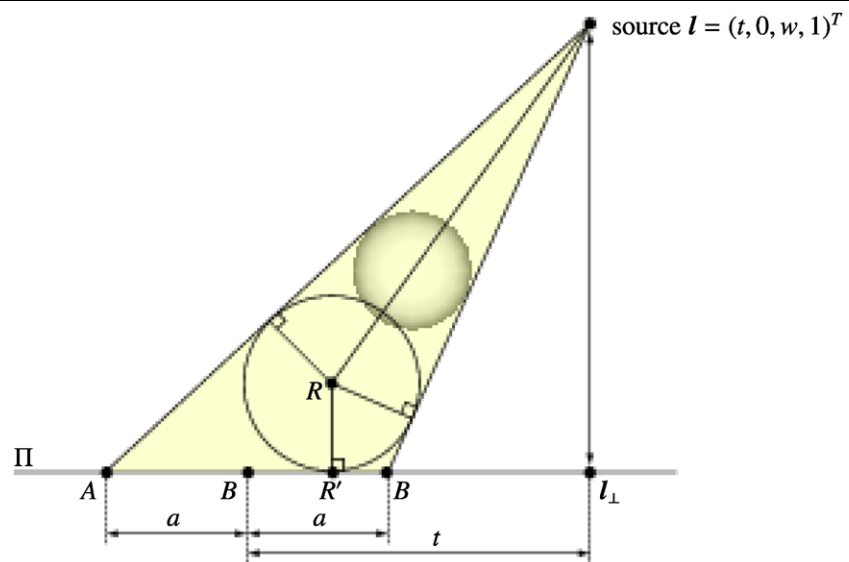


Fig. 7 Source positions (u, w) are estimated using three calibration spheres. The sizes and positions of the spheres and screen are shown in the plot. Each plot shows 11 source positions obtained from (a) ground truth, (b) accurate shadowgrams, and (c)–(d) shadowgrams blurred using 5×5 and 10×10 averaging filters. On the right is the visual hull of a branch reconstructed from 50 light sources. The poor result demonstrates the need for better algorithms for reconstructing intricate shapes

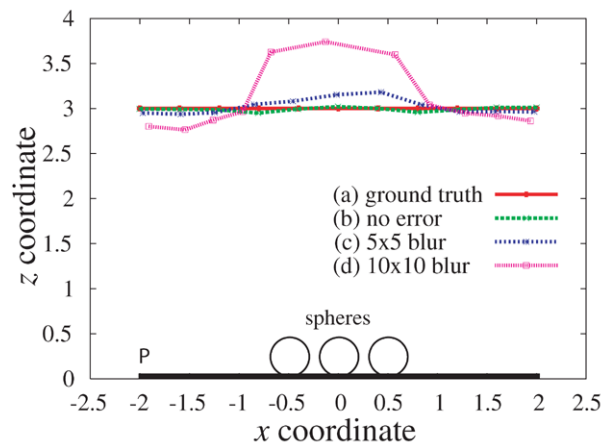


Fig. 8 Epipolar geometry of two shadowgrams. The baseline connecting the two sources l_i and l_j intersects the shadowgram plane Π at an epipole e_{ij} . Suppose an epipolar plane that is tangent to the surface of an object at a frontier point F , then the intersection of the epipolar plane and the shadowgram plane Π is an epipolar line. The epipolar line can be estimated as a line that is co-tangent to the shadowgrams at T_i and T_j

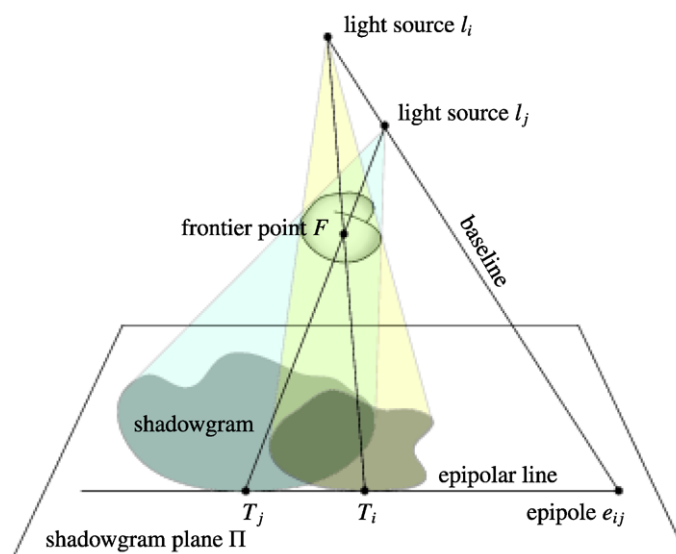


Fig. 9 Localization of the epipole. (a), (b) If two shadowgrams are convex, a maximum of four co-tangent lines and six intersections are possible. Considering that the object and the light source are on the same side with respect to the screen, the epipole can be chosen uniquely out of the six intersections. (c), (d) If the shadowgrams are non-convex, the epipole is localized by applying the technique in (a) or (b) to the convex polygons of the original shadowgrams

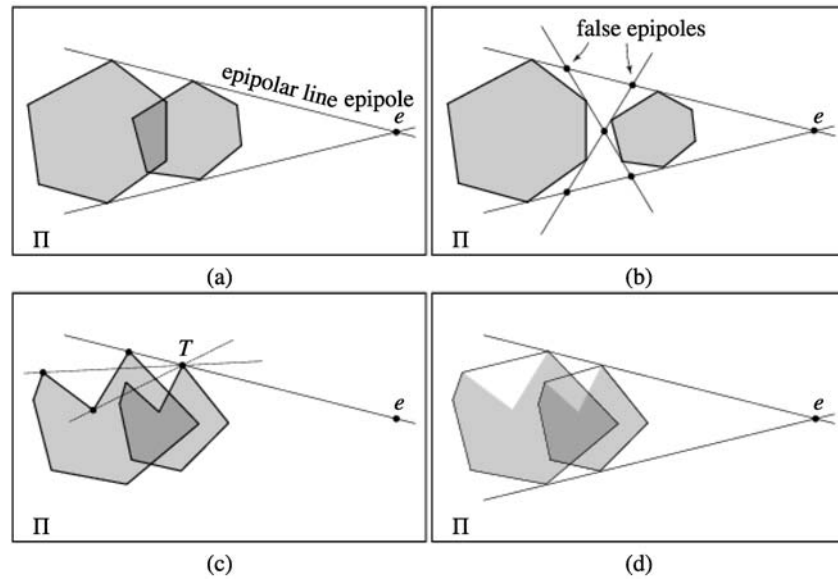


Fig. 10 Initial light source positions in Fig. 7 were improved by epipolar constraints in (24). On the right is the visual hull reconstructed from the improved source positions

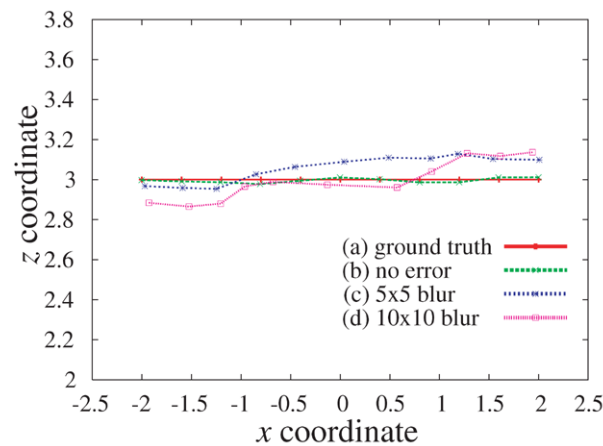
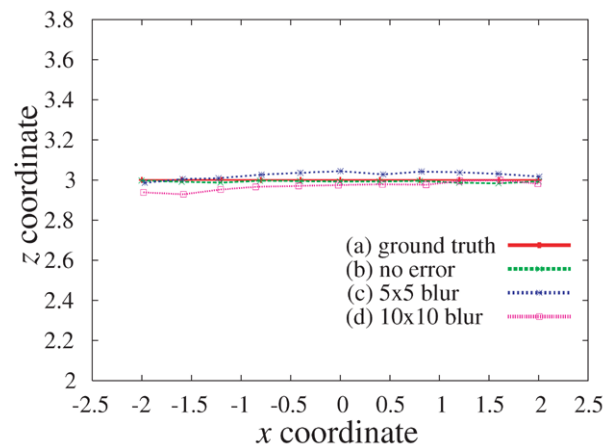


Fig. 11 The light source positions reconstructed using epipolar constraint in Fig. 10 were optimized by maximizing the shadowgram consistency in (28). On the right is the visual hull reconstructed from the optimized source positions



this case, there are no frontier points formed (and hence no cotangent lines for convex silhouettes). We can avoid this case by placing the sources such that the baselines do not always intersect the object.

4.2 Improving Accuracy of Source Locations

The error in the light source positions reconstructed using spheres can be arbitrarily large depending on the localiza-

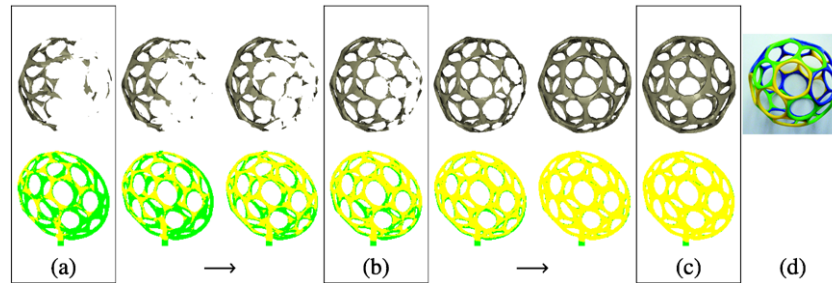
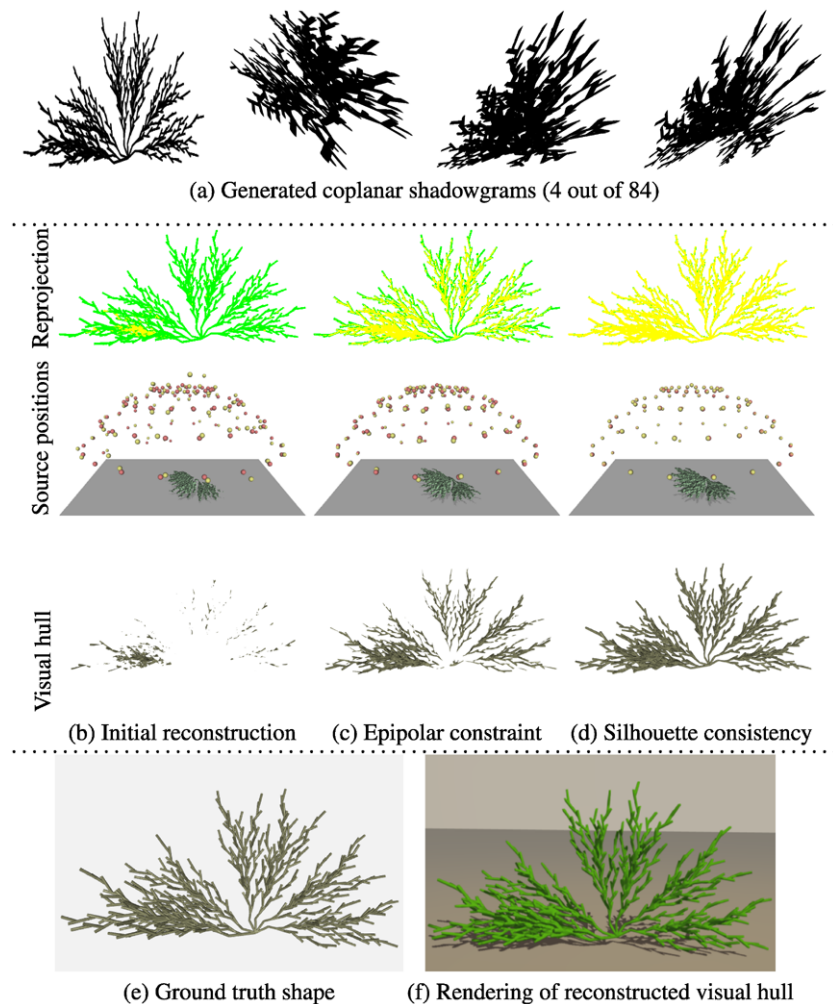


Fig. 12 Reconstructed shape of a polyhedron object is improved with each iteration from left to right. (Top) Reconstructed visual hull at the end of each iteration. (Bottom) The reprojection of the reconstructed visual hulls onto one of captured silhouette images. The reprojection and silhouettes are consistent at yellow pixels, and inconsistent at

green (Color online). The boxed figures show the reconstruction from the light source positions (a) estimated from spheres, (b) improved by epipolar geometry, and (c) optimized by maximizing shadowgram consistency. (d) Photograph of the object

Fig. 13 Simulation with a coral object: (a) Eighty four coplanar shadowgrams of the object are generated with average resolution 530×270 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) The ground truth 3D shape. The volume difference is 0.15% of the volume of the ground truth 3D shape. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)



tion of the elliptical shadowgram for each sphere. This error can be reduced by relating different light source positions through the epipolar geometry. Let the set of epipoles e_{ij} be estimated from all the source pairs l_i and l_j . The locations of the sources are improved by minimizing the expression in

(19) for each pair of light sources using least squares:

$$\{l_1, \dots, l_N\} = \underset{\{l_1, \dots, l_N\}}{\operatorname{argmin}} \sum_{i \neq j} \|M_{ij} e_{ij}\|_2^2 \quad (24)$$

where $\|\cdot\|_2$ is the L2-norm of a vector. The source positions reconstructed from the shadowgrams of spheres are used as

Fig. 14 Simulation with a seaweed object: (a) Forty nine coplanar shadowgrams of the object are generated with average resolution 334×417 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) The ground truth 3D shape. The volume difference is 0.21% of the volume of the ground truth 3D shape. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)

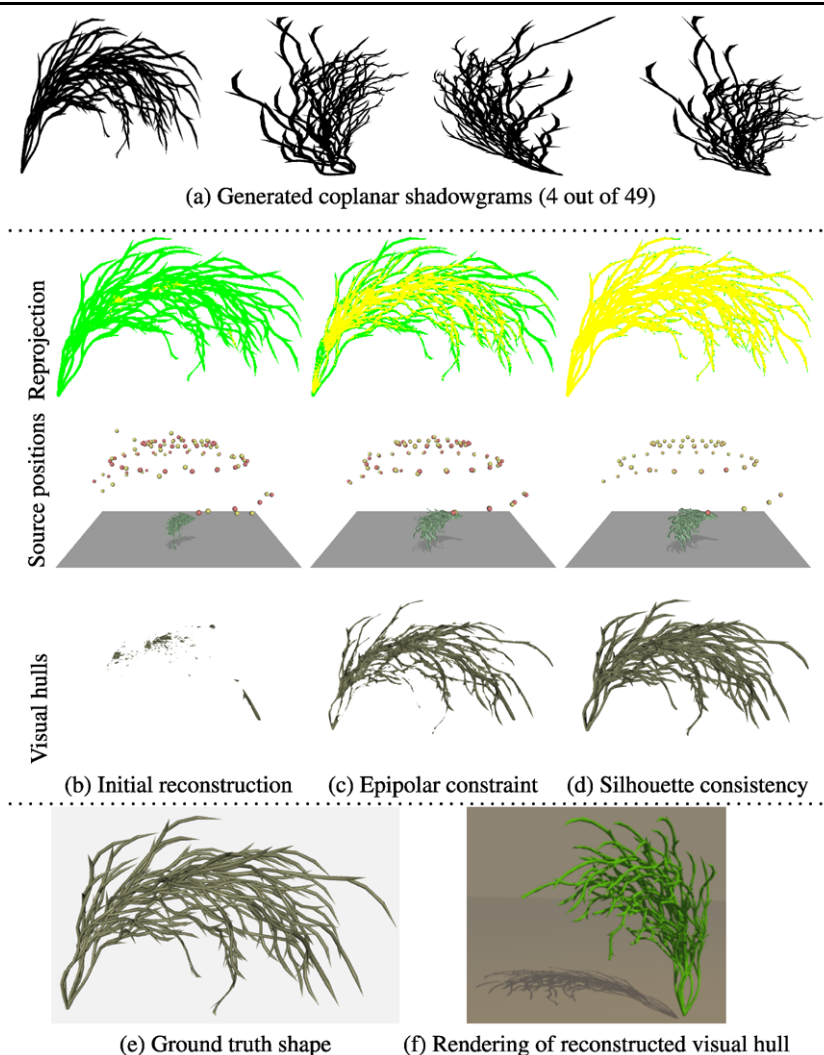


Table 2 Comparison between traditional multi-view camera-based imaging and coplanar shadowgrams in epipolar reconstruction. The traditional multi-view images acquired by uncalibrated cameras re-

quire at least 7 point correspondences of the silhouette contours. Coplanar shadowgrams allow unique epipolar reconstruction irrespective of the shape of the 3D object

Silhouette complexity #correspondences	Convex	Non-convex		
	2	< 7	≥ 7	$\gg 7$
Uncalibrated multi-camera	impossible	impossible	not always	hard
Coplanar shadowgrams	possible	possible	possible	possible

possible—the epipolar geometry can be reconstructed uniquely

not always—possible if seven correspondences are found

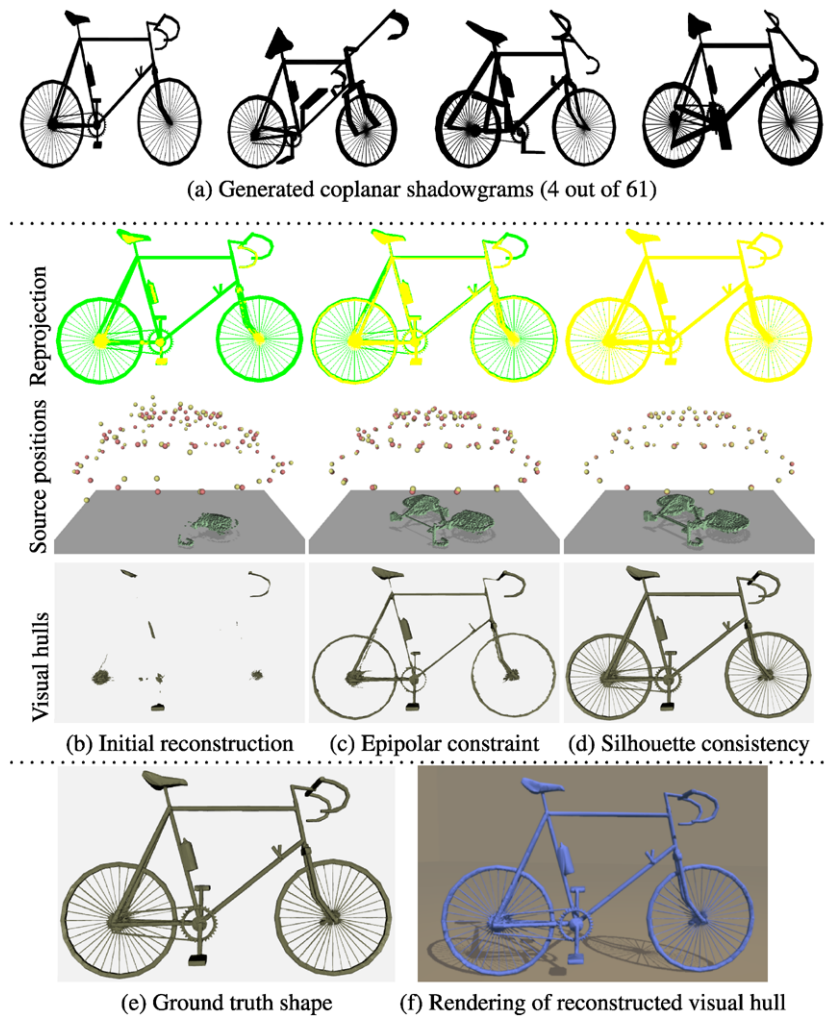
hard—hard to find the correct correspondences in practice

impossible—impossible because of the insufficient constraints

initial estimates. We evaluate this approach using the simulated silhouettes described in Fig. 7. Figure 10 shows considerable improvement in accuracy obtained by enforcing

the epipolar constraint in (19). Compared to the result in Fig. 7, collinearity in the positions of light sources is better recovered in this example.

Fig. 15 Simulation with a **bicycle** object: **(a)** Sixty one coplanar shadowgrams of the object are generated with average resolution 635×425 pixels. **(b)** Initial reconstruction. **(c)** The reconstruction using epipolar geometry. **(d)** The reconstruction using silhouette consistency. **(e)** The ground truth 3D shape. The volume difference is 0.12% of the volume of the ground truth 3D shape. **(f)** Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)



5 Using Shadowgram Consistency

While the epipolar geometry improves the estimation of the light source positions, the accuracy of estimate can still be insufficient for the reconstruction of intricate shapes (Fig. 10). In this section, we present an optimization algorithm that improves the accuracy of all the source positions even more significantly. As we will show, combining the epipolar constraints and the optimization algorithm results in high quality models of intricate shapes.

5.1 Optimizing Light Source Positions

Let \mathcal{V} be the visual-hull obtained from the set of captured shadowgrams \mathcal{S}_i and the estimated projection $P(l_i)$ for $i = 1, \dots, N$. Due to the nature of the intersection operator, the reprojections of the visual hull to the shadowgram plane satisfy:

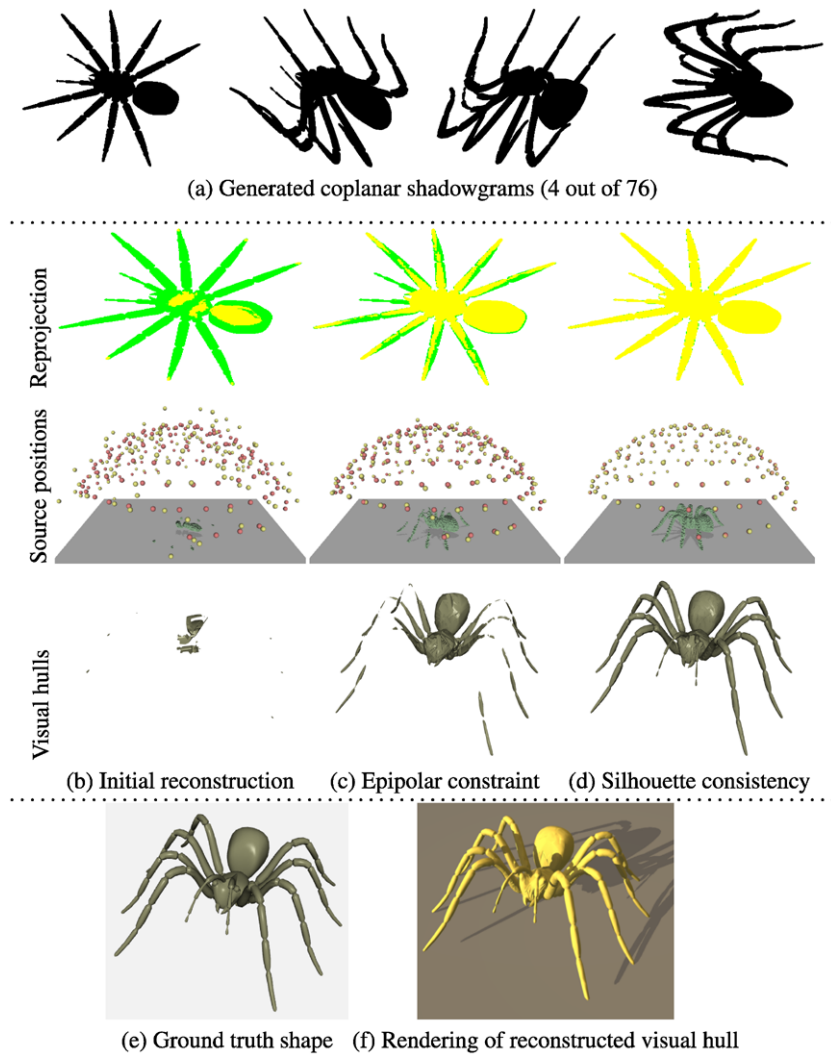
$$\mathcal{S}_i^{\mathcal{V}} \subseteq \mathcal{S}_i^{\mathcal{O}}. \quad (25)$$

The reprojections match the acquired silhouettes when the source positions are perfect. In other words, if they match, we cannot make any more improvement to the shape only by the silhouettes. We can define a measure of silhouette mismatch by the sum of squared difference:

$$\epsilon_{\text{reprojection}}^2 = \sum_{i=1}^N D(\mathcal{S}_i^{\mathcal{V}}(x), \mathcal{S}_i^{\mathcal{O}}(x)) \quad (26)$$

where $D: \mathbb{RP}^2 \times \mathbb{RP}^2 \rightarrow \mathbb{R}$ evaluates the difference between two sets and will be discussed in Sect. 5.2. We minimize the above mismatch by optimizing for the locations of the light sources. Unfortunately, optimizing solely (26) is known to be inherently ambiguous owing to the four-parameter transformation mentioned in Sect. 3. To alleviate this issue, we simultaneously minimize the discrepancy between the optimized light source positions l_i and the initial source positions l_i^0 estimated from the spheres (Sect. 3) and

Fig. 16 Simulation with a spider object: (a) Seventy six coplanar shadowgrams of the object are generated with average resolution 356×354 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) The ground truth 3D shape. The volume difference is 0.08% of the volume of the ground truth 3D shape. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)



epipolar geometry (Sect. 4):

$$\epsilon_{\text{initial}}^2 = \sum_{i=1}^N \|l_i - l_i^0\|_2^2. \quad (27)$$

The final objective function is obtained by a linear combination of the two errors:

$$\epsilon_{\text{total}} = \epsilon_{\text{reprojection}}^2 + \alpha \epsilon_{\text{initial}}^2 \quad (28)$$

where α is a user-defined weight. While the idea of minimizing silhouette discrepancy is well known in the traditional multi-view camera-based SFS (Sinha et al. 2004; Yezzi and Soatto 2003; Wong and Cipolla 2004; Hernández et al. 2007), the key advantage over prior work is the reduced number of parameters our algorithm needs to optimize (three per view for the light source position, instead of six per view for rotation and translation of the camera). In turn, this allows us to apply our technique to a much larger number of views than possible before.

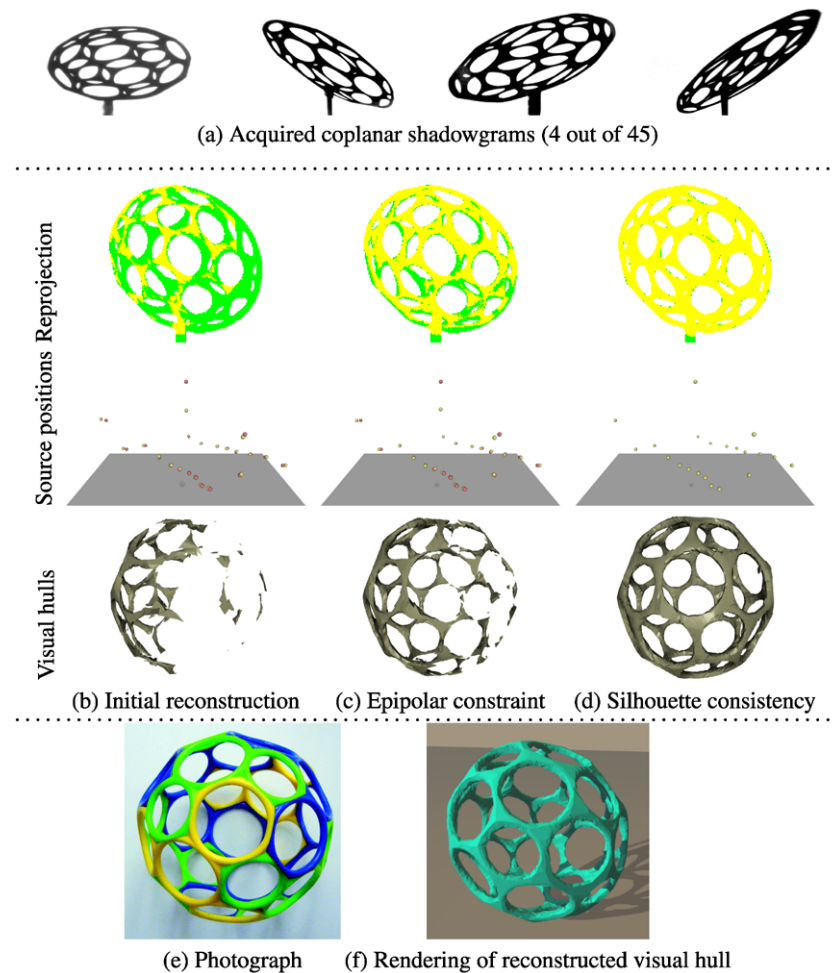
5.2 Implementation

How do we represent the scalar-valued function D that describes the *difference* between two silhouettes in (26)? Suppose the silhouettes \mathcal{S}_i^V and \mathcal{S}_i^O are given as images in the same dimension, then it is natural to define D as the sum of square distances between the pixels values of the silhouettes. Then, the next question is how to define the pixels values of the silhouette images.

The simplest one is a binary function that returns 0 and 1 when a pixel is located respectively outside and inside the silhouette. This binary function, however, is not suitable for iterative minimization of (26) since $D(\mathcal{S}_i^V, \mathcal{S}_i^O)$ is non-zero at only a small number of pixels around mis-matched contour lines. Due to errors in the measured silhouette, the binary function causes undesirable bias.

We use the signed Euclidean distances to the contour of \mathcal{S}_i^V and \mathcal{S}_i^O as the pixels values in the silhouette images. The intersection of silhouettes is computed for each 3D ray defined by a pixel in \mathcal{S}_i^O , and then projected back to

Fig. 17 Real experiment with a polyhedron object: (a) Forty five coplanar shadowgrams of the object are generated with average resolution 126×116 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) Photograph of the object. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)



the silhouette to obtain S_i^V . This is a simplified version of image-based visual hull (Matusik et al. 2000) and has been used in silhouette registration methods (Hernández et al. 2007). Equation (28) is minimized using Powell's gradient-free technique (Press et al. 1988).

Due to the intricate shapes of the silhouettes, the error function in (28) can be complex and may have numerous local minima. We alleviate this issue using the convex polygons of the silhouette contours described in Sect. 4. From Proposition 1, a following corollary is derived regarding the consistency between the visual hull and the silhouettes of an object (see Appendix D for the proof).

Proposition 2 *If silhouette contours are consistent in that they can be generated from a physical 3D object, then the convex polygons obtained from the silhouette contours are also consistent.*

Using Proposition 2, we minimize (28) using the convex silhouettes with I_i^0 as initial parameters. The resulting light source positions are in turn used as starting values to min-

imize (28) with the original silhouettes. In practice, using convex silhouettes also speeds up convergence.

We evaluate this approach using the simulated silhouettes described in Figs. 7 and 10. Compare the results in Fig. 7 (using spheres to estimate source positions) and Fig. 10 (enforcing epipolar constraints) with those in Fig. 11. The final reconstruction of the tree branch is visually accurate highlighting the performance for our technique.

6 Results

In this section, we demonstrate the accuracy of our techniques using both simulated and real experimental data. Table 3 summarizes the data set used in the experiment and the performance of our reconstruction algorithms.

We first generated or captured a sufficiently large number of shadowgrams, and started our reconstruction algorithm using a small number of randomly-chosen images. The number of images used is increased until the reconstructed visual gets sufficiently close to the actual shape. We used a workstation equipped with four dual-core AMD Opteron 8218

Fig. 18 Real experiment with a wreath object: (a) 122 coplanar shadowgrams of the object are generated with average resolution 674×490 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) Photograph of the object. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)

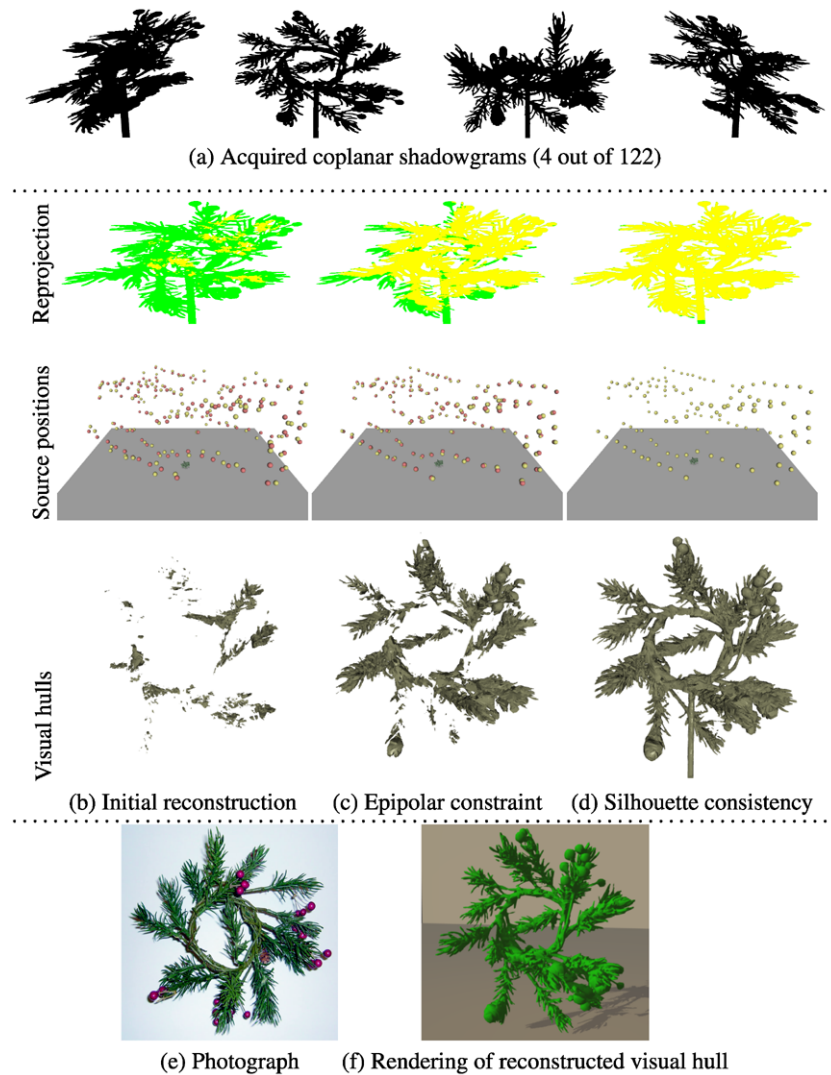


Table 3 The models used in our experiments: We reconstructed visual hulls from coplanar shadowgrams using four simulation and four real data sets. The detail of each experiment is shown in the corresponding figure. The average resolution of shadow region in the shadowgrams is shown in the row of shadowgram size. The computation time in optimizing epipolar geometry and silhouette consistency is shown in

minutes (see the main text for our computational environment). Reprojection error indicates the mismatch between the input shadowgrams and those generated by reprojecting the estimated visual hull. For simulation data, the volume ratio between the ground truth and the reconstructed visual hulls is presented









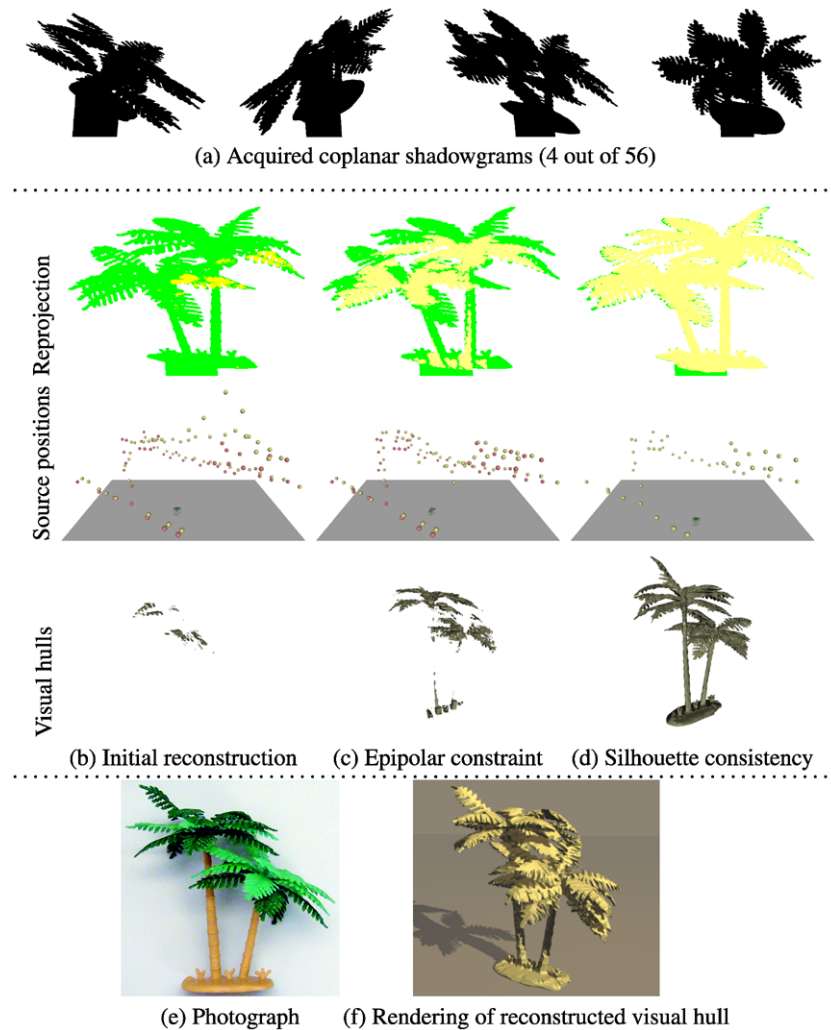
	Simulation data				Real data			
Model	coral	seaweed	bicycle	spider	polyhedron	wreath	palm-tree	octopus
								
Figure	13	14	15	16	17	18	19	20
# of views	84	49	61	76	45	122	56	53
Shadowgram size	530×270	334×417	635×425	356×354	126×116	674×490	520×425	451×389
Reconstruction time	532 min	384 min	429 min	498 min	235 min	9312 min	429 min	402 min
Reprojection error	2.2%	3.2%	2.3%	1.3%	3.2%	5.2%	4.8%	4.6%
Volume error	0.15%	0.21%	0.12%	0.08%	–	–	–	–

Fig. 19 Real experiment with a palm-tree object: (a) Fifty six coplanar shadowgrams of the object are generated with average resolution 520×425 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) Photograph of the object. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)



2.6 GHz and 16 GB main memory for optimization of epipolar geometry and silhouette consistency explained in Sect. 4 and Sect. 5.

All results of 3D shape reconstructions shown in this paper are generated by the exact polyhedral visual hull method proposed by Franco and Boyer (2003). The acquired 3D shape is then rendered using Autodesk Maya rendering package.

6.1 Reconstruction of Visual Hulls and 3D Source Positions

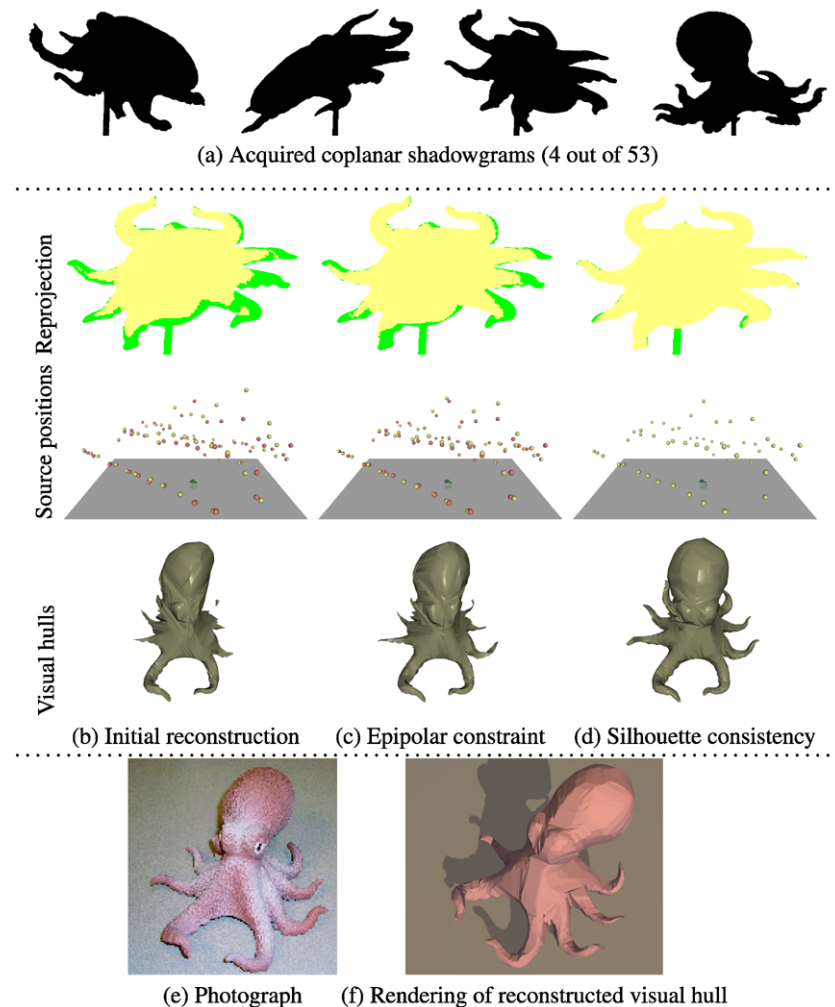
6.1.1 Simulation Data

We have chosen four objects with complex structure in our simulations—a coral, a seaweed (also used in Figs. 7, 10, and 11 in the main paper), a bicycle, and a spider. The seaweed and coral objects have many thin sub-branches with numerous occlusions. The bicycle object is composed of very thin structures such as spokes, chains, and gears. The

spider object is composed of both thick and thin structure. The simulation experiments with known ground truth shape and source positions are shown respectively in Figs. 13, 14, 15, and 16.

Each of the figures is organized as follows: (a) A set of coplanar shadowgrams of the object is generated by a shadow simulator implemented by Direct 3D graphics library. (b) The positions of light source are perturbed with random noise with $\sigma = 5\%$ of the object size, and the silhouettes are blurred by 3×3 averaging filters. (c) The positions of the light sources are recovered using epipolar geometry followed by the maximization of silhouette consistency in (d). For each of (b), (c), and (d), the top row shows one of captured silhouette images (in green), overlaid with the reprojection of the reconstructed visual hulls onto the silhouette (in yellow). The middle row shows the ground truth positions of light sources (in red) and the estimated positions (in yellow). The reconstructed 3D shape is shown at the bottom. Finally, (e) the ground truth 3D shape and (f) the reconstructed visual hull rendered by Maya is shown.

Fig. 20 Real experiment with an octopus object: (a) Fifty three coplanar shadowgrams of the object are generated with average resolution 451×389 pixels. (b) Initial reconstruction. (c) The reconstruction using epipolar geometry. (d) The reconstruction using silhouette consistency. (e) Photograph of the object. (f) Rendering of the reconstructed shape. (Refer to main text for the detail of each figure)



6.1.2 Real Data

We show the 3D shape reconstruction of four different objects—a polyhedron (also used in the Fig. 12 in the main paper), a wreath (Figs. 1 and 2), a palm-tree object, and an octopus object. The wreath object has numerous thin needles which cause severe occlusions. The polyhedron is a thin wiry polyhedral object. The palm-tree object is a plastic object composed of two palm trees with flat leaves. The octopus is a plastic object that has complex surface reflection and large concavities. The results of reconstructing 3D shape are shown in Figs. 17, 18, 19, and 20. Each figure is organized in the same way as those of simulation data, except that: The final reconstruction of source positions are presented in red in the middle row of (b), (c), and (d). The photograph of the object is shown in (e).

6.2 Convergence

Figure 12 illustrates the convergence properties of our optimization algorithm. Figure 12(a) shows the visual hull of the

wiry polyhedral object obtained using the initial positions of light sources estimated from the calibration spheres. The reprojection of the visual hull shows poor and incomplete reconstruction. By optimizing the light source positions, the quality of the visual hull is noticeably improved in only a few iterations.

The convergence of the reconstruction algorithm is quantitatively evaluated in Fig. 21. The error in light source positions estimated by the algorithm proposed in Sect. 5 is shown in the top-left plot. The vertical axis shows L2 distance between the ground truth and the current estimate of light source positions. After convergence, the errors in the light source positions are less than 1% of the sizes of the objects. The volume ratio between actual and erroneously-reconstructed visual hulls is presented in the top-right. The silhouette mismatch defined in (26) is plotted on the bottom. On average, the silhouettes cover on the order of 10^5 pixels. The error in the reprojection of the reconstructed visual hulls is less than 1% of the silhouette pixels for the real objects.

Fig. 21 Convergence of error: (a, b) Errors in light source positions and visual hull are computed using ground truth for simulation models. (a) The error in the source position computed by L2 distance in meter. (b) The volume ratio between actual and erroneously-reconstructed visual hulls. (c) Error in shadowgram consistency computed in pixels for both simulation and real data. All plots are in logarithmic scale

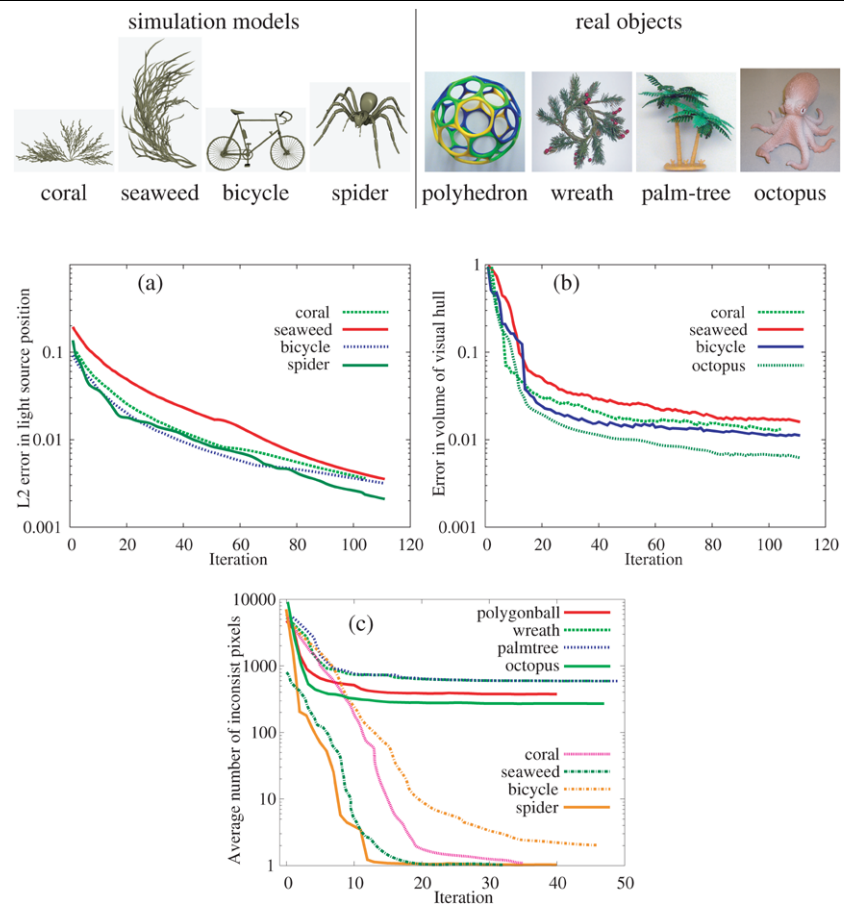
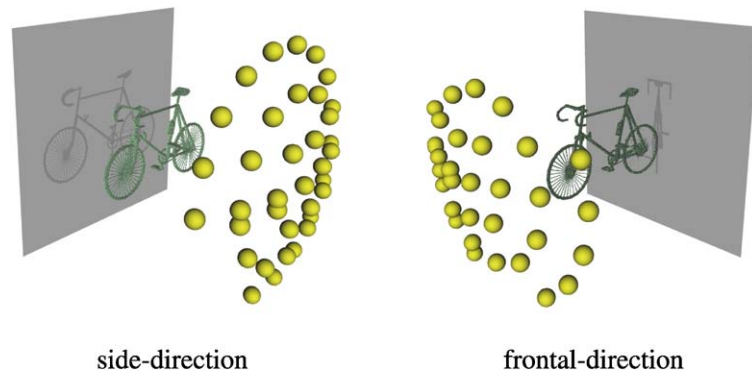


Fig. 22 Two different configurations of coplanar shadowgrams of a bicycle object: Gray rectangle and yellow spheres indicate respectively a shadow screen and light source position (Color online). 36 light sources are used in both configurations. The screen is rotated by 90 degrees, while the object remains fixed. For the demonstration of the two-screen algorithm, a small number of light sources are used



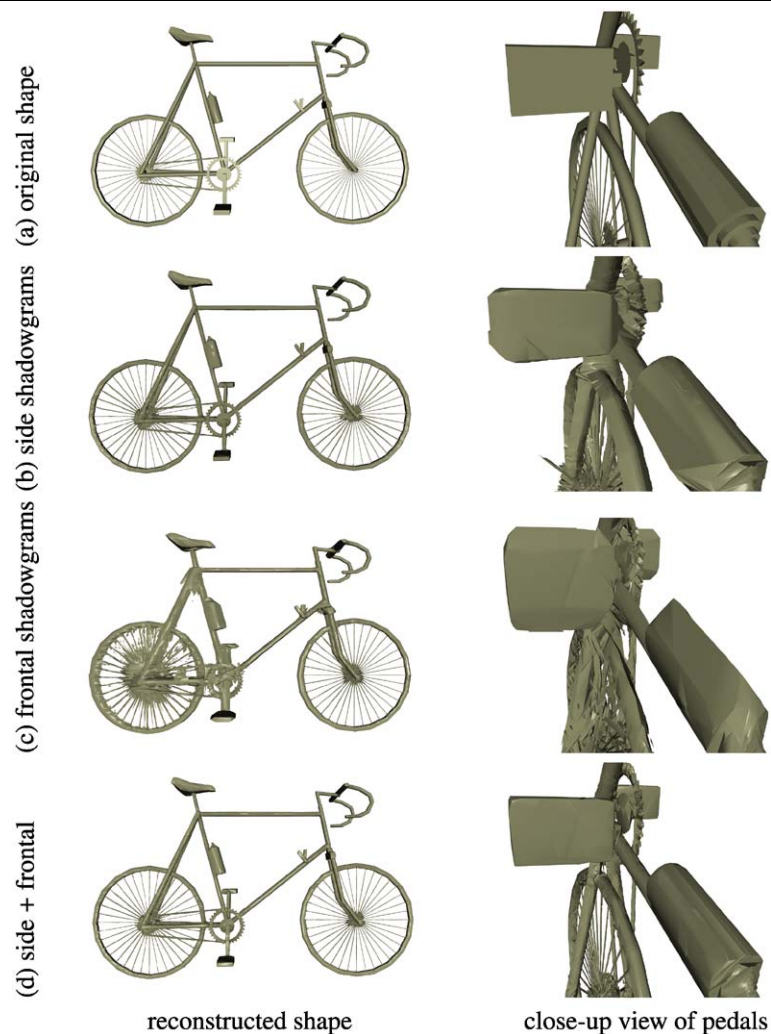
7 Conclusion

7.1 Summary

Coplanar shadowgram imaging is a technique to easily capture a large number of silhouettes of intricate objects from a wide range of viewpoints, with high accuracy. The setup shown in Fig. 4 is inexpensive and requires only off-the-shelf components like a camera, a rear-projection screen and a small light source. However, traditional shape-from-silhouettes (SFS) approaches are highly sensitive to (even tiny) errors in acquired silhouettes or light source posi-

tions. The epipolar geometry of the shadowgrams and the silhouette consistency based optimization described in this work are crucial to obtaining extremely accurate reconstructions of intricate shapes. Most approaches for shape-from-silhouettes have acquired shapes of non-intricate objects (simple models of people, statues and figurines) that can be modeled with a small number of views. We believe this is the first time such accurate shapes have been estimated automatically using a large number of views (50 to 120) of objects with severe occlusions, discontinuities and thin structures.

Fig. 23 Comparison of shape reconstruction: We synthesized 36 coplanar shadowgrams of a 3D shape shown in (a). The visual hull of the object is reconstructed from: (b) the shadowgrams taken from side-direction (Fig. 22 left) and (c) the shadowgrams taken from frontal-direction (Fig. 22 right). The reconstructed shape is stretched into the direction perpendicular to a shadow screen due to the lack of views parallel to the screen. (d) Combining shadowgrams (b) and (c) enlarges the coverage of light source positions, which successfully reduces the stretching artifact in the reconstructed shape



7.2 Discussion

A single screen cannot be used to capture the complete $360^\circ \times 360^\circ$ view of the object. For instance, it is not possible to capture the silhouette observed in the direction parallel to a shadowgram plane. This limitation can be overcome by augmenting the system with more than one shadowgram screen (or move one screen to different locations). The algorithm of the multi-screen coplanar shadowgram imaging can be divided into offline and online steps:

Off-line Calibration (one-time): This calibration can be done in several ways and we mention a simple one here. In the case of two-screen setup which is observed by a single camera, we only need to estimate the homography between each screen and image plane. The extra work required over the one-screen case is an additional homography estimation. The homographies in turn can be used to recover the relative transformation between the screens.

Online Calibration: In the two-screen setup, we can estimate the light source positions for each set of shadowgrams

on one screen separately using the technique demonstrated in the paper. Finally, we merge the two sets of results using the relative orientation between the screens resulting from the off-line calibration.

In principle, it is possible to also optimize (minimize) the errors due to off-line calibration. However, the off-line intrinsic calibration of a camera and the screen-to-image homography can be done carefully. More importantly, it is independent of the complexity of the object and the number of source positions.

We have performed simulations with a bicycle object with two screen positions as shown in Fig. 22. The bicycle was chosen since frontal and side views are both necessary to carve the visual hull satisfactorily. Combining the two sets of shadowgrams enlarges the coverage of source positions, which successfully reduces the stretching artifact of the reconstructed shape in Fig. 23.

7.3 Future Work

One drawback of SFS techniques is the inability to model concavities on the object's surface. Combining our approach with other techniques, such as photometric stereo or multi-view stereo can overcome this limitation, allowing us to obtain appearance together with a smoother shape of the object. Finally, using multiple light sources of different spectra to speed up acquisition, and the analysis of penumbra due to the finite size of a light source are our directions of future work.

Acknowledgements This research was supported by ONR award N00014-05-1-0188 and ONR DURIP award N00014-06-1-0762. Narasimhan is also supported by NSF CAREER award IIS-0643628. Kanade is supported by NSF grant EEE-540865.

Appendix A: Derivation of (2)

The projective transformation in coplanar shadowgram imaging is viewed as a perspective transformation in the translated coordinate system whose origin is at the location of a light source. Thus, the projection by the source located at $\mathbf{l} = (u, v, w, 1)^T$ is written in matrix form as

$$\begin{aligned} P(\mathbf{l}) &= \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}}_{\text{drop } z \text{ } (\equiv 0)} \cdot \underbrace{\begin{pmatrix} \mathbf{I}_3 & \mathbf{l} \\ \mathbf{0}_3^T & 1 \end{pmatrix}}_{\text{project to } \Pi} \cdot \begin{pmatrix} -w\mathbf{I}_3 & \mathbf{0}_3 \\ (0, 0, 1) & 0 \end{pmatrix} \\ &\quad \cdot \begin{pmatrix} \mathbf{I}_3 & -\mathbf{l} \\ \mathbf{0}_3^T & 1 \end{pmatrix} \\ &= \begin{pmatrix} -w & 0 & u & 0 \\ 0 & -w & v & 0 \\ 0 & 0 & 1 & -w \end{pmatrix} \end{aligned}$$

where \mathbf{I}_3 is a 3×3 identity matrix and $\mathbf{0}_3 = (0, 0, 0)^T$.

Appendix B: Derivation of (12)

The linear transformation \mathbf{A} is generally written in a 4×4 matrix as

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{pmatrix}. \quad (29)$$

Then, the elements in (11) can be expanded into a 3×4 matrix using (2) and (29), which yields twelve equations. Using the six equations that are identical to 0 in the left side, we obtain

$$a_{21} = a_{31} = a_{41} = a_{12} = a_{32} = a_{42} = a_{14} = a_{24} = a_{34} = 0.$$

Similarly,

$$a_{11} = a_{22} = a_{44} \neq 0 \quad (30)$$

because $w \neq 0$ and $w' \neq 0$. Dividing \mathbf{A} by $a_{11} (= a_{22} = a_{44} \neq 0)$, we can reparameterize \mathbf{A} to (12). By solving (11) with respect to \mathbf{l}' , the transformation from \mathbf{l} to \mathbf{l}' can also be obtained as (14). To guarantee that \mathbf{l} and \mathbf{l}' are in the same side with respect to the shadowgram plane, the additional constraints on the elements in \mathbf{A} is derived as (13).

Appendix C: Proof of Proposition 1

Proof

Suppose the silhouette $\mathcal{S}_i^{\mathcal{O}}$ of an object \mathcal{O} is generated by perspective projection P_i ($i = 1, \dots, N$). The convex hull of the object and the silhouettes are respectively defined as

$$\hat{\mathcal{O}} \stackrel{\text{def}}{=} \bigcup_{x_j, x_k \in \mathcal{O}} \overline{x_j x_k} \quad \text{and} \quad (32)$$

$$\hat{\mathcal{S}}_i^{\mathcal{O}} \stackrel{\text{def}}{=} \bigcup_{m_j, m_k \in \mathcal{S}_i^{\mathcal{O}}} \overline{m_j m_k} \quad (33)$$

where the over-line $\overline{p_j p_k}$ represents the line segment spanned by two points p_j and p_k . Then, the shadowgram $\mathcal{S}_i^{\hat{\mathcal{O}}}$ of the convex object $\hat{\mathcal{O}}$ generated by the projection P_i is

$$\mathcal{S}_i^{\hat{\mathcal{O}}} \stackrel{\text{def}}{=} P_i \cdot \hat{\mathcal{O}} \quad (34)$$

$$= \bigcup_{x_j, x_k \in \mathcal{O}} P_i \cdot \overline{x_j x_k} \quad (\because (32)) \quad (35)$$

$$= \bigcup_{m_j, m_k \in \mathcal{S}_i^{\mathcal{O}}} \overline{m_j m_k} \quad (\because P_i \text{ is linear.}) \quad (36)$$

$$= \hat{\mathcal{S}}_i^{\mathcal{O}} \quad (\because (33)). \quad (37)$$

□

Appendix D: Proof of Proposition 2

Proof Suppose the visual hull \mathcal{V} of an object \mathcal{O} is reconstructed from the shadowgrams $\mathcal{S}_i^{\mathcal{O}}$ and the projections P_i ($i = 1, \dots, N$). Proposition 1 implies that the convex hull $\hat{\mathcal{O}}$ of the object generates the convex hulls $\hat{\mathcal{S}}_i^{\mathcal{O}}$ of the shadowgrams:

$$\hat{\mathcal{S}}_i^{\mathcal{O}} = P_i \cdot \hat{\mathcal{O}}. \quad (38)$$

Hence, the convex shadowgrams $\hat{\mathcal{S}}_i^{\mathcal{O}}$ are consistent by definition. □

References

- Åström, K., Cipolla, R., & Giblin, P. (1999). Generalised epipolar constraints. *International Journal of Computer Vision*, 33(1), 51–72.
- Balan, A., Sigal, L., Black, M., & Haussecker, H. (2007). Shining a light on human pose: On shadows, shading and the estimation of pose and shape. In *Proc. international conference on computer vision '07*.
- Baumgart, B. G. (1974). *Geometric modeling for computer vision*. PhD thesis, Stanford University.
- Besant, W. H. (1890). *Conic sections, treated geometrically*. Cambridge: Deighton, Bell.
- Boykov, Y., & Funka-Lea, G. (2006). Graph cuts and efficient n-d image segmentation. *International Journal of Computer Vision*, 70(2), 109–131.
- Campbell, N., Vogiatzis, G., Hernández, C., & Cipolla, R. (2007). Automatic 3d object segmentation in multiple views using volumetric graph-cuts. In *Proc. the British machine vision conference '07* (pp. 530–539).
- Chandraker, M. K., Kahl, F., & Kriegman, D. J. (2005). Reflections on the generalized bas-relief ambiguity. In *Proc. computer vision and pattern recognition '05* (Vol. 1, pp. 788–795).
- Cheung, K. M., Baker, S., & Kanade, T. (2005). Shape-from-silhouette across time, Part I: Theory and algorithms. *International Journal of Computer Vision*, 62(3), 221–247.
- Cipolla, R., Åström, K., & Giblin, P. (1995). Motion from the frontier of curved surfaces. In *Proc. international conference on computer vision '95* (pp. 269–275).
- Cross, G., Fitzgibbon, A. W., & Zisserman, A. (1999). Parallax geometry of smooth surfaces in multiple views. In *Proc. international conference on computer vision '99* (pp. 323–329).
- Curless, B., & Levoy, M. (1996). A volumetric method for building complex models from range images. In *Proc. SIGGRAPH '96* (pp. 303–312).
- Drbohlav, O., & Chantler, M. (2005). Can two specular pixels calibrate photometric stereo? In *Proc. international conference on computer vision '05* (pp. 1850–1857).
- Drbohlav, O., & Sara, R. (2002). Specularities reduce ambiguity of uncalibrated photometric stereo. In *Proc. the 7th European conference on computer vision* (pp. 46–62).
- Fitzgibbon, A., Pilu, M., & Fisher, R. (1999). Direct least squares fitting of ellipses. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(5), 476–480.
- Franco, J.-S., & Boyer, E. (2003). Exact polyhedral visual hulls. In *Proc. the 15th British machine vision conference* (pp. 329–338).
- Furukawa, Y., Sethi, A., Ponce, J., & Kriegman, D. (2006). Robust structure and motion from outlines of smooth curved surfaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(2), 302–315.
- Georghiades, A. S. (2003). Incorporating the torrance and sparrow model of reflectance in uncalibrated photometric stereo. In *Proc. international conference on computer vision '03* (Vol. 2, pp. 816–823).
- Hartley, R., & Zisserman, A. (2004). *Multiple view geometry in computer vision*, 2nd edn. Cambridge: Cambridge University Press.
- Hayakawa, H. (1994). Photometric stereo under a light-source with arbitrary motion. *JOSA*, 11(11), 3079–3089.
- Hernández, C., Schmitt, F., & Cipolla, R. (2007). Silhouette coherence for camera calibration under circular motion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(2), 343–349.
- Hooke, R. (1667). *Micrographia*. London, Chap. Observation LVIII.
- Kriegman, D. J., & Belhumeur, P. N. (2001). What shadows reveal about object structure. *Journal of the Optical Society of America*, 18(8), 1804–1813.
- Levoy, M., Pulli, K., Curless, B., Rusinkiewicz, S., Koller, D., Pereira, L., Ginzton, M., Anderson, S., Davis, J., Ginsberg, J., Shade, J., & Fulk, D. (2000). The digital michelangelo project: 3D scanning of large statues. In *Proc. SIGGRAPH '00* (pp. 131–144).
- Matusik, W., Buehler, C., Raskar, R., Gortler, S. J., & McMillan, L. (2000). Image-based visual hulls. In *Proc. SIGGRAPH '00* (pp. 369–374).
- Press, W., Flannery, B., Teukolsky, S., & Vetterling, W. (1988). *Numerical recipes in C*. Cambridge: Cambridge University Press.
- Savarese, S., Andreetto, M., Rushmeier, H., Bernardini, F., & Perona, P. (2005). 3D reconstruction by shadow carving: Theory and practical evaluation. *International Journal of Computer Vision*, 71(3), 305–336.
- Sawhney, H. S. (1994). Simplifying motion and structure analysis using planar parallax and image warping. In *Proc. international conference of pattern recognition* (pp. 403–408).
- Seitz, S. M., Curless, B., Diebel, J., Scharstein, D., & Szeliski, R. (2006). A comparison and evaluation of multi-view stereo reconstruction algorithms. In *Proc. computer vision and pattern recognition '06* (Vol. 1, pp. 519–526).
- Settles, G. S. (2001). *Schlieren & shadowgraph techniques*. Berlin: Springer.
- Sinha, S. N., Pollefeys, M., & McMillan, L. (2004). Camera network calibration from dynamic silhouettes. In *Proc. computer vision and pattern recognition '04* (Vol. 1, pp. 195–202).
- Smith, A. R., & Blinn, J. F. (1996). Blue screen matting. In *Proc. SIGGRAPH '96* (pp. 259–268).
- Tan, P., Mallick, S., Kriegman, D., Quan, L., & Zickler, T. (2007). Isotropy, reciprocity and the gbr ambiguity. In *Proc. computer vision and pattern recognition '07* (pp. 1–8).
- Wong, K.-Y. K., & Cipolla, R. (2004). Reconstruction of sculpture from its profiles with unknown camera positions. *IEEE Transactions on Image Processing*, 13(3), 381–389.
- Yamazaki, S., Narasimhan, S., Baker, S., & Kanade, T. (2007). Coplanar shadowgrams for acquiring visual hulls of intricate objects. In *Proc. international conference on computer vision '07*, October 2007.
- Yezzi, A. J., & Soatto, S. (2003). Stereoscopic segmentation. *International Journal of Computer Vision*, 1(53), 31–43.