

TOWARDS STRUCTURE AND MOTION ESTIMATION FROM DYNAMIC SILHOUETTES *

Tanuja Joshi Narendra Ahuja Jean Ponce
Beckman Institute,
University Of Illinois, Urbana, IL 61801

Abstract: *We address the problem of estimating the structure and motion of a smooth curved object from its silhouettes observed over time by a trinocular imagery. We first construct a model for the local structure along the silhouette for each frame in the temporal sequence. The local models are then integrated into a global surface description by estimating the motion between successive frames. The algorithm tracks certain surface and image features (parabolic points and silhouette inflections, frontier points) which are used to bootstrap the motion estimation process. The whole silhouette is then used to refine the initial motion estimate. We have implemented the proposed approach and report preliminary results.*

1 Introduction

Structure and motion estimation for objects with smooth surfaces and little texture is an important but difficult problem. Silhouettes are the dominant image features of such objects. The viewing cone grazes the smooth surface along the (occluding) contour and intersects the image plane along the silhouette. At each contour point, the surface normal is perpendicular to the viewing direction. If the object is moving relative to the viewer, the silhouettes appearing in successive images will be projections of different 3D contours on the surface. This is in contrast with an image sequence containing only viewpoint-independent features.

Several methods have been proposed for structure estimation from silhouettes under *known* camera motion [1, 4, 12, 13, 14]. These approaches have demonstrated that given a set of three or more nearby views of a smooth object, the structure of the object up to second order can be obtained along its silhouette. The recovery of structure *and* motion from a monocular sequence of silhouettes has been investigated by Giblin et al. [3]. For the case of a curved object rotating about a fixed axis with constant angular velocity, they show: (1) given a complete set of orthographic silhouettes, the rotation axis and velocity can be recovered, along with the visible surface; (2) given the silhouettes over a short time interval, the rotation axis can be recovered if the angular velocity is known.

*This work was supported in part by the Advanced Research Projects Agency under Grant N00014-93-1-1167 and in part by the National Science Foundation under Grant IRI-9224815.

In this paper, we address the problem of estimating the structure and motion of a smooth object undergoing *arbitrary* unknown rigid motion from its silhouettes observed over time by a trinocular stereo rig. This technique will be useful when the viewer is a passive observer and has no knowledge or control of the object's motion (see [10] for a complementary approach, where a viewer plans its motion for building a global model of an object). Another application is model construction for object recognition[7]: due to self-occlusion, a simple motion, such as rotation on a turntable, may not reveal all the interesting parts of a complex object; it is desirable to move the object arbitrarily and still be able to construct the complete model.

We use trinocular imagery for our analysis since three images can be used to recover the local structure up to second order. The world coordinate frame is taken to be that of the central camera of the trinocular imagery. We assume orthographic projection. At a given time t , the three images taken by the trinocular imagery are used to estimate the local structure of the object along the silhouette in the central image. If the motion of the object is small, the structure estimated at time t will also be valid at time $t + 1$. We use the structure computed at time t to estimate the motion between time frames t and $t + 1$.

Strictly speaking, since the 3D contour changes with the object's relative motion, it is impossible to define a unique point-to-point correspondence between successive silhouettes. Like others [1, 12, 13], we use the epipolar curves to establish correspondences and estimate the local structure. When the relative motion is unknown, we rely instead on other surface curves that project onto detectable and trackable silhouette features to obtain an initial estimate of the unknown motion; we then iteratively refine the motion estimate using the rest of the silhouette.

The silhouette *inflections* form such a set of features. They are projections of parabolic points on the surface [8, 13]. We use the change in the surface normals at the matched inflections to obtain an estimate of the rotation between two successive time frames. Given the epipolar plane geometry for a pair of images, we can consider another set of features called *frontier points*, where the surface normal is parallel to the epipolar plane normal [3].

In Sect. 2 the algorithm for structure estimation using trinocular imagery is described. The algorithm used for motion estimation from dynamic silhouettes is discussed in Sect. 3. To demonstrate the feasibility of our methods, we

include experimental results on a set of synthetic images. We conclude with comments in Sect. 4.

2 Structure Estimation Using Trinocular Imagery

2.1 Modeling the Local Structure

The local structure (up to second order) at a surface point P is defined by the 3D location of P , the surface normal at P , the two principal directions in the tangent plane and the principal curvatures at P . At each point P , we define a local coordinate frame (x_l, y_l, z_l) with its origin at P , the x_l -axis along the outward surface normal, and the y_l - and z_l -axes along the two principal directions. The local surface up to second order is a paraboloid [2], given by $x_l = (\kappa_1 y_l^2 + \kappa_2 z_l^2)/2$ (where κ_1 and κ_2 are the principal curvatures at P) or in matrix form given by:

$$\mathbf{Q}_l^T \mathbf{M}_l \mathbf{Q}_l = 0, \quad (1)$$

where \mathbf{M}_l is a symmetric 4×4 matrix and \mathbf{Q}_l is the vector of homogeneous coordinates of a point Q on the paraboloid at P .¹ The signs of κ_1 and κ_2 define the point type: if κ_1 and κ_2 have the same sign (resp. opposite signs), P is an elliptic (resp. hyperbolic) point. If either κ_1 or κ_2 is zero, P is a parabolic point and the silhouette has an inflection.

The rigid transformation parameters between the local and the world coordinate frames (in our case the camera-centered frame), together with the principal curvatures κ_1 and κ_2 , completely describe the local structure at P . Let (x, y, z) be the camera-centered frame, with the z -axis along the viewing direction and the xy -plane being the image plane. In general we need six parameters for the rigid transformation between (x_l, y_l, z_l) and (x, y, z) frames. But if P is a contour point, by definition the surface normal (the x_l -axis) is orthogonal to the viewing direction (the z -axis). Hence we need only five parameters: two rotational and three translational.

Consider again a point P on the contour (see Fig. 1). Let the angle between the x_l - and x -axes be θ and the angle between the z_l - and z -axes be γ . Let P be at (x_0, y_0, z_0) in the camera-centered frame. The five-tuple $(\theta, \gamma, x_0, y_0, z_0)$ defines the local frame with respect to the camera-centered frame. To completely describe the surface locally, we need to specify κ_1 and κ_2 in addition to the above five-tuple for each point on the silhouette. Equation 1 can be rewritten in (x, y, z) frame as:

$$\mathbf{Q}^T \mathbf{M} \mathbf{Q} = 0, \quad (2)$$

where $\mathbf{M} = \mathbf{T}_0^{-1T} \mathbf{R}_z \mathbf{R}_\theta \mathbf{M}_l \mathbf{R}_\theta^{-1} \mathbf{R}_z^{-1} \mathbf{T}_0^{-1}$.

Here \mathbf{R}_θ is the 4×4 matrix in homogeneous coordinates for a rotation of angle γ about the x_l -axis, \mathbf{R}_z is the 4×4

¹Notation: All boldface letters denote coordinate vectors or arrays. All capital letters denote 3D points in the scene and small letters denote their projections in the image.

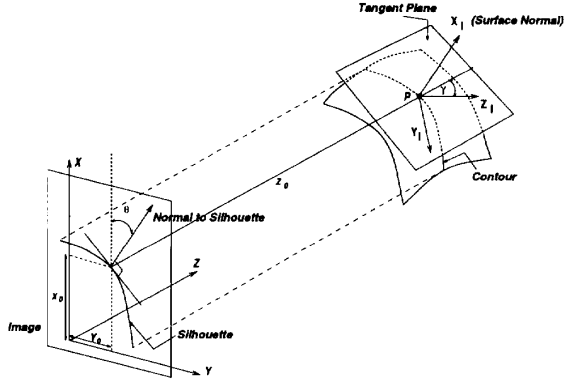


Figure 1: Projection geometry.

matrix for a rotation of angle θ about the z_l -axis, \mathbf{T}_0 is the 4×4 matrix for a translation by (x_0, y_0, z_0) , and \mathbf{Q} is the homogeneous coordinate vector of Q in the camera-centered frame.

Since P is a contour point, the surface normal at P is orthogonal to the viewing direction (z -axis), hence the z -component of the surface normal is zero:

$$\frac{d(\mathbf{Q}^T \mathbf{M} \mathbf{Q})}{dz} = 0. \quad (3)$$

This is a linear condition in (x, y, z) , implying that the contour of the paraboloid is a planar curve. Eliminating z between (2) and (3) gives the equation of the silhouette, which is a parabola [6]. From this single projection in the central frame, we can obtain three constraints on the seven structure parameters $(\theta, \gamma, x_0, y_0, z_0, \kappa_1, \kappa_2)$: the normal to the silhouette at the apex of the parabola is parallel to the surface normal, which gives angle θ directly. The image location of the apex gives x_0 and y_0 . The curvature at the apex gives a constraint on κ_1, κ_2 and γ . To complete the local structure model, we need to estimate depth z_0 and obtain two more constraints on κ_1, κ_2 , and γ which are obtained using the matched points in the other two images of the trinocular imagery.

2.2 Finding Correspondences

Let I_1, I_2 and I_3 be the three images taken by the three cameras of the trinocular imagery. Since the relative positions of the cameras are known, we can define an epipolar plane for each point for each pair of images. Consider a silhouette point p_1 in image I_1 . Similar to conventional stereo, we can find the correspondence p_2 (resp. p_3) in image I_2 (resp. I_3) as the point lying on the epipolar line corresponding to p_1 . The difference here is that the matched image points are not projections of the same 3D point.

For the case of a continuous relative motion between a camera and an object, at each point on the silhouette the

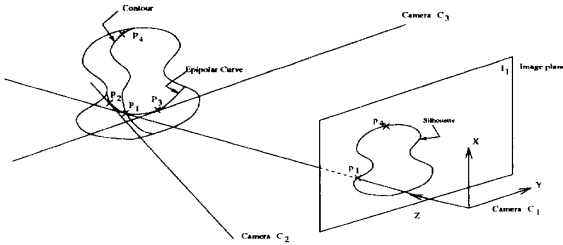


Figure 2: Epipolar geometry.

points matched using the instantaneous epipolar geometry trace a 3D curve (called epipolar curve) on the object such that at each point the curve is tangent to the view line [1, 12, 13]. In the trinocular imagery case, the points matched using the epipolar geometry will lie on the corresponding epipolar curve (see Fig. 2). At a point with the surface normal not parallel to the epipolar plane normal this curve is a regular curve e.g. at point P_1 in Fig. 2. But at points where the surface normal is parallel to the epipolar plane normal, the epipolar curve degenerates into a point e.g. at point P_4 in Fig. 2. In fact, in such a case, the matched points are projections of the same 3D point. Such points are called the *frontier points* for the corresponding camera motion [3].

As we approach a frontier point, the tangent to the silhouette becomes parallel to the epipolar line, making the computation of the epipolar match difficult and inaccurate. Here we use the paraboloid model to find the correspondence. It can be shown that the projection of the local paraboloid is a parabola in each image with its normal at the apex depending only on the surface normal and the relative rotational motion of the camera [6]. Given a frontier point in image I_1 and the relative position and orientation of camera C_2 (resp. C_3) we can predict the normal to the corresponding parabola in image I_2 (resp. I_3). We find the matching parabola (in effect the matching point) in I_2 (resp. I_3) using the predicted normal.

2.3 Computing Structure Parameters

At Non-frontier points: Previous approaches for estimation of structure under known camera motion [1, 12, 13] have used the epipolar parameterization of the surface: in [1] a differential formulation is presented to construct the structure along the silhouette; in [13] the radial plane is used to estimate one of the surface curvatures.

Our reconstruction method is similar to the one used by Szeliski and Weiss [12], where the epipolar plane is chosen instead of the radial plane for the computation of one of the surface curvatures. In our setup, the three viewing directions are taken to be coplanar for simplicity, giving a common epipolar plane for the three cameras and making the epipolar curve planar. This is not an essential assumption. If the view lines are not coplanar, they can

be projected onto the epipolar plane, but this involves an approximation.

Since the epipolar curve is tangent to the three view lines, we can estimate its osculating circle by finding the circle that is tangent to the three view lines. The point where this circle touches the central view line is an estimate of the 3D surface point. This gives depth z_0 of the 3D point. The curvature of the fitted circle is an estimate of the curvature of the epipolar curve. This curvature enables us to compute the normal curvature of the surface along the view line, which in turn gives a constraint on the surface curvatures κ_1, κ_2 and angle γ from the Euler formula [2].

Once the depth of the points along the silhouette is computed, we can estimate the direction of the tangent to the 3D contour. The contour tangent gives one more constraint on κ_1, κ_2 and γ since it is along a direction conjugate to the view line [9]. This constraint, along with the constraints given by the normal curvature and the curvature of the silhouette in the central image give us three equations which are solved to obtain the values of the structural parameters, κ_1, κ_2 and γ . See [6] for further details.

At Frontier Points: We make use of the paraboloid model described in Sect. 2.1 to estimate the structure at a frontier point. It can be shown [6] that the projection of the local paraboloid in images I_2 and I_3 is also a parabola. For each of these parabolas: (1) The normal at the apex is given by the surface normal of the parabola in I_1 and the relative rotational motion between the two cameras. As discussed in Sec. 2.2, we use this fact in matching the parabolas (in effect points). (2) The curvature at the apex gives a constraint on the surface curvatures. The three constraints obtained from the three matched parabolas are solved to compute κ_1, κ_2 and γ . (3) The location of the apex gives a constraint on depth z_0 . This enables us to obtain two constraints for depth z_0 from the locations of the matches in the other two images. We solve these individually and then take the average of the two computed values to get the final estimate of z_0 . See [6] for further details.

2.4 Experimental Results

In this preliminary implementation, we have applied the method described in this section to a set of synthetic images generated using two fourth degree algebraic surfaces shown in Fig. 3 a-b. For a given viewing geometry, the 3D contour of the object was computed using the numerical method described in [11]. The computed contour was then projected on a 512×512 image under orthographic projection. The images, although synthetic, contained quantization noise. The edge points in the images were linked to give a closed curve and a local cubic model was fit at each point to estimate the tangent and curvature of the silhouette. The inflections were detected as the zero-crossings of the estimated curvature.

Figure 3.c-d shows sample reconstructed contours along

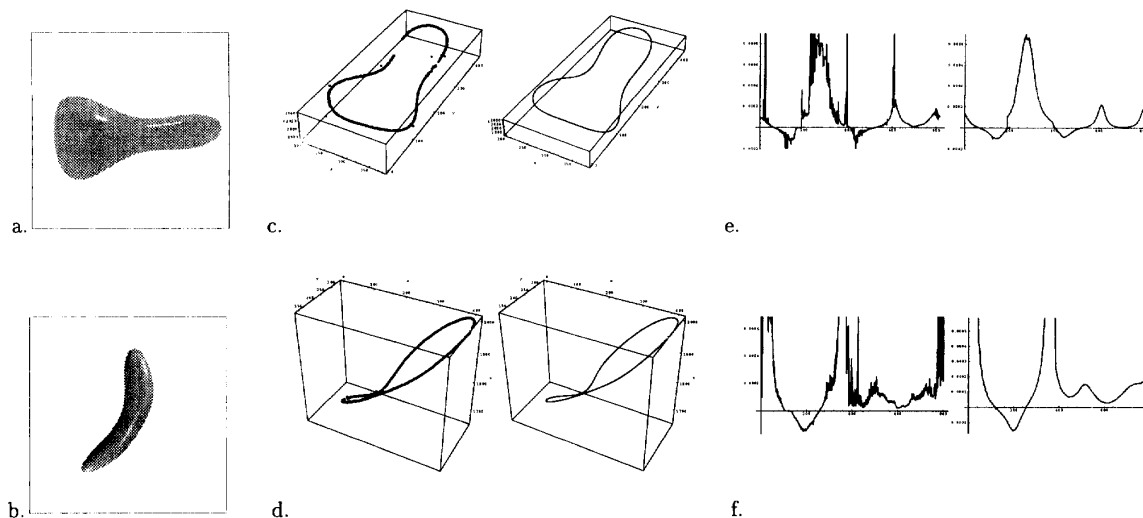


Figure 3: Results of structure estimation: a-b. two objects, c-d. sample recovered and true 3D contours for the two objects, e-f. recovered and true Gaussian curvatures along the two contours

with the true contours. The gaps in the reconstructed contours correspond to parabolic points where the surface normal is nearly parallel to the epipolar normal. It is not possible to get the depth estimate at these points (see [5] for a discussion on epipolar curves and where the reconstruction based on epipolar parameterization fails). Fig. 3.e-f shows the estimated and true Gaussian curvatures along the silhouette.

3 Motion Estimation

Let's assume that between consecutive time frames, the object has undergone a rotation of an unknown angle α about an unknown axis $\Omega = [\omega_x, \omega_y, \omega_z]^T$ passing through the origin, followed by an unknown translation $[t_x, t_y, t_z]^T$. Let R and T be the 4×4 matrices corresponding to the unknown rotation and translation.

We have 3D contours on the surface in the successive frames $I_1(t)$ and $I_1(t+1)$ of the central camera estimated using the methods in Sect. 2. To estimate the relative motion between successive frames, we use two sets of features of the silhouette: the sets of inflections and frontier points.

The motion estimation is done in two steps. In the first step, we use the change in the surface normal at the inflections to estimate the rotation parameters. For a given rotation, the epipolar plane normal can be estimated, and in turn the frontier points can be detected. These are used to estimate the translation parameters. In the second step, the estimate of the motion is refined using the entire silhouette. Both the steps perform non-linear least-squares minimization over the rotation parameters. The next two

sections describe these two steps in detail. The experimental results are reported in Sect. 3.3.

3.1 Obtaining Initial Motion Estimate

Inflections are projections of parabolic points (points with zero Gaussian curvature) on the object surface. On generic surfaces these points lie on continuous curves - called parabolic curves. A parabolic point has a single asymptotic direction (which is also one of the principal directions) and the surface is locally cylindrical.

We use the fact that at a parabolic point, if we move on the surface along any direction (hence in particular, along the direction tangent to the parabolic curve), the change in the surface normal is orthogonal to the asymptotic direction at that point [6, 7].

Consider an inflection p in the central image $I_1(t)$ which is the projection of a parabolic point P . If we track the inflection in image $I_1(t+1)$, we will be moving along the parabolic curve at P . From the above stated fact, the change in the normal should be orthogonal to the asymptotic direction A at P . We can compute A from the local structure at P estimated at time t . Let p' be the tracked inflection in $I_1(t+1)$, and the projection of a surface point P' with the unit surface normal N' . Thus N' has to satisfy the constraint the de-rotated vector $R^{-1}N'$ is perpendicular to A . Here R is a matrix representing the unknown rotation of the object between time t and $t+1$.

We parameterize the rotation using three parameters - angle ϕ between the rotation axis and the z -axis, angle ψ between its projection in xy -plane and the x -axis, and

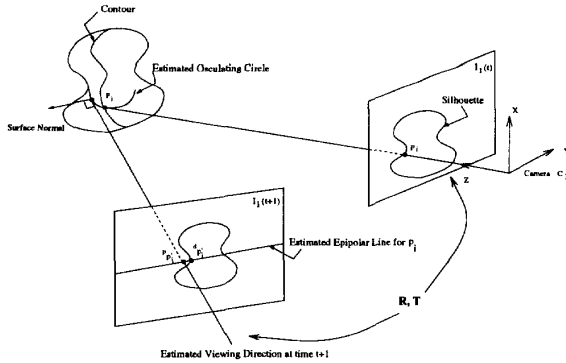


Figure 4: Predicted and detected epipolar match.

the rotation angle α . To recover the rotation parameters, we need to track at least three inflections. With $n \geq 3$ inflections present in images $I_1(t)$ and $I_1(t+1)$, we use least-squares minimization with the objective function:

$$\sum_{i=1}^n [A_i \cdot (R^{-1} N'_i)]^2. \quad (4)$$

The minimization is done over the three parameters ϕ , ψ and α . We have to first find a match between the sets of inflections on the silhouettes in the two images. We match the inflections such that (1) their ordering along the silhouette is maintained, and (2) the angle between the normals at the matched inflections is small.

For a given rotation, we can estimate the translation using the frontier points of the silhouette. The corresponding frontier points are projections of the same 3D point on the surface. Hence if a frontier point f in image $I_1(t)$ matches with the frontier point f' in image $I_1(t+1)$, we have $F' = T R F$. Thus for a given rotation, estimating the translation becomes a linear problem once the matches for the frontier points are known. We match the frontier points such that their order along the silhouette is maintained and the residual error in estimating the translation is minimum.

3.2 Refining the Motion Estimate

In the second step of motion estimation, we refine the initial estimate using the structure along the entire silhouette. With the complete estimate of R and T , we can determine the viewing direction of frame $I_1(t+1)$ relative to the local paraboloids in frame $I_1(t)$, and hence in turn determine the epipolar plane for each point in image $I_1(t)$. Knowing the structure parameters of the local paraboloids we can estimate the curvature of the epipolar curve at each point as well.

Consider a point p_i on the silhouette in image $I_1(t)$ (see Fig. 4 where we have shown the object-centered view for convenience). We can predict the match point ${}^p p'_i$ in frame

$I_1(t+1)$ as the projection of the point which (1) lies on the estimated osculating circle of the epipolar curve and (2) has the surface normal orthogonal to the viewing direction of $I_1(t+1)$. But we can also detect the epipolar match point ${}^d p'_i$ in image $I_1(t+1)$ using the estimated epipolar geometry. In the refinement step, we iteratively minimize the sum of the squared distance between ${}^p p'_i$ and ${}^d p'_i$ for each silhouette point p_i . Thus the algorithm for refining the estimate can be stated as follows.

1. For the given rotation parameters ϕ , ψ and α , detect the frontier points f and f' in the two frames. Find the match for the set of frontier points. Using this match, compute the translation T .
2. Using the estimate of R and T and the local structure at each point p_i on the silhouette, predict the epipolar match point ${}^p p'_i$ at $t+1$. Also detect the epipolar match point ${}^d p'_i$ in image $I_1(t+1)$. Compute the sum, say s , of the squared distance between ${}^p p'_i$ and ${}^d p'_i$ for all the silhouette points p_i .
3. Repeat steps 1 and 2 to minimize s , updating the values of angles ϕ , ψ and α at each step.

3.3 Experimental Results

We have tried a variety of motions to test our algorithm. It was found that the first step of minimization was stable with respect to the initial guess. The result of this step is used as the initial guess for the second step. We noticed that the estimate of the rotation axis after the first step is within 15 degrees of the correct axis. Therefore we repeat the second minimization step with the initial guess for the rotation axis sampled evenly from a cone of 15 degrees around the result of the first step. The result giving minimum residual error was taken to be the final result.

Table 1 lists the recovered motion parameters after each step on a sample set of motions. For each step, we also list the angle between the true and estimated rotation axes and the error in the rotation angle. The estimate of angle α was within one degree of the correct value. Since the rotation is modeled to be about an axis passing through the origin, a small error in the rotation parameters results into a relatively large error in translational parameters. Figure 5.a-c shows the recovered 3D contours de-rotated after the motion estimate, and overlaid on the true 3D contours for the set of motions given in Table 1. The recovered contours fit closely to the true ones.

We also applied the algorithm of motion estimation over a sequence of five frames. Fig. 5.d shows all the detected contours (along with the recovered local paraboloids) in the object-center frame de-rotated using the estimated motion. The small errors in the interframe motion estimation have not caused accumulation of errors over frames.

4 Conclusions

Although estimating motion from silhouettes is more difficult than using viewpoint-independent features, we

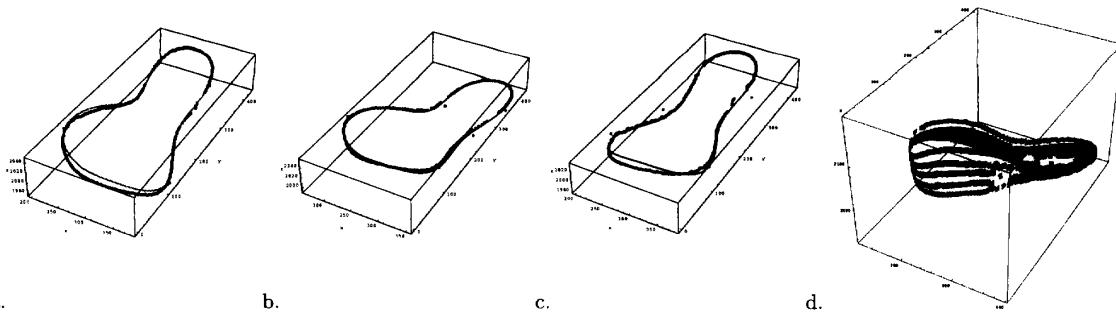


Figure 5: Results of motion estimation: a-c. de-rotated contours overlaid on the true contours, d. global structure estimation using a sequence of five frames.

		Rotation Axis $\Omega(\omega_x, \omega_y, \omega_z)$	Error in Ω	Rotation Angle α	Error in α	Translation (t_x, t_y, t_z) in pixels
1	True	(0.58, 0.58, 0.58)	-	-10.00	-	(190.38, -210.64, 20.26)
	1st step	(0.61, 0.69, 0.39)	12.68	-13.05	3.05	(288.40, -289.03, 56.04)
	2nd step	(0.68, 0.53, 0.51)	7.53	-10.99	0.99	(178.76, -270.66, 35.67)
2	True	(0.58, 0.58, 0.58)	-	10.00	-	(-210.64, 190.38, 20.266)
	1st step	(0.31, 0.72, 0.62)	17.67	8.78	1.22	(-231.40, 89.75, 13.82)
	2nd step	(0.60, 0.58, 0.54)	2.46	10.55	0.55	(-228.27, 215.81, 18.83)
3	True	(-0.00, 0.71, 0.71)	-	10.00	-	(-245.57, -15.19, 15.19)
	1st step	(-0.02, 0.74, 0.67)	3.06	11.18	1.18	(-270.30, -24.73, 5.49)
	2nd step	(0.06, 0.73, 0.68)	3.99	10.83	0.83	(-266.48, 5.53, 2.42)

Table 1: Result of the motion estimation (All angles are in degrees)

have the advantage that even from a single silhouette we can obtain more information about the surface: the surface normal, the sign of the Gaussian curvature and a constraint on the principal curvatures at each contour point. Our method uses the dynamic silhouettes to estimate both the motion and the global surface of a smooth object.

The initial experimental results are encouraging and provide a demonstration of the validity of the method. More rigorous experiments should be conducted on synthetic images with controlled noise and on real images. Since three frames are sufficient to distinguish between viewpoint-dependent and viewpoint-independent edges [13], using the trinocular imagery we can detect viewpoint-independent edges present, if any, and utilize them in motion estimation as well. We also plan to investigate the case of perspective projection.

References

- [1] A. Blake and R. Cipolla. Surface shape from the deformation of apparent contours. *Int. J. of Comp. Vision*, 9(2), 1992.
- [2] M. P. do Carmo. *Differential Geometry of Curves and Surfaces*. Prentice-Hall, Englewood Cliffs, NJ, 1976.
- [3] P. Giblin, J. Rycroft, and F. Pollock. Moving surfaces. In *Mathematics of Surfaces V*. Cambridge University Press, 1993.
- [4] P. Giblin and R. Weiss. Reconstruction of surfaces from profiles. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, 1987.
- [5] P. Giblin and R. Weiss. Epipolar fields on surfaces. In *Proc. European Conf. Comp. Vision*, 1994.
- [6] T. Joshi, N. Ahuja, and J. Ponce. Structure and motion estimation from dynamic silhouettes. Technical Report UIUC-BI-AI-RCV-94-01, Beckman Institute, University of Illinois, 1994.
- [7] T. Joshi, J. Ponce, B. Vijayakumar, and D. Kriegman. Hot curves for modelling and recognition of smooth curved 3d shapes. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, 1994.
- [8] J. J. Koenderink. What does the occluding contour tell us about solid shape. *Perception*, 13, 1984.
- [9] J. J. Koenderink. *Solid Shape*. MIT Press, MA, 1990.
- [10] K. Kutulakos and C. Dyer. Global surface reconstruction by purposive control of observer motion. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, 1994.
- [11] S. Petitjean, J. Ponce, and D.J. Kriegman. Computing exact aspect graphs of curved objects: Algebraic surfaces. *Int. J. of Comp. Vision*, 9(3), 1992.
- [12] R. Szeliski and R. Weiss. Robust shape recovery from occluding contours using a linear smoother. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, 1993.
- [13] R. Vaillant and O. Faugeras. Using extremal boundaries for 3-d object modeling. *IEEE Trans. Patt. Anal. Mach. Intell.*, 14(2), 1992.
- [14] J. Y. Zheng. Acquiring 3-d models from sequences of contours. *IEEE Trans. Patt. Anal. Mach. Intell.*, 16(2), 1994.