# Model-Based Brightness Constraints: On Direct Estimation of Structure and Motion

Gideon P. Stein and Amnon Shashua, Member, IEEE

Abstract—We describe a new direct method for estimating structure and motion from image intensities of multiple views. We extend the direct methods of Horn and Weldon [18] to three views. Adding the third view enables us to solve for motion and compute a dense depth map of the scene, directly from image spatio-temporal derivatives in a linear manner without first having to find point correspondences or compute optical flow. We describe the advantages and limitations of this method which are then verified with experiments using real images

Index Terms—Shape representation and recovery, 3D recovery from 2D, shape from motion, image sequence analysis, algebraic and projective geometry.

# **1** INTRODUCTION

**T**HE geometry of multiple views of a 3D scene is wellunderstood. There exist geometric constraints, which relate corresponding features (points and lines) in multiple views to the camera geometry. These constraints take the form of the trilinear tensor equations for three views and the epipolar constraints for two views. Given a set of feature correspondences, we can recover the camera geometry and, in particular, the camera motion (or the relative position of multiple cameras). But finding correspondences is a hard problem, and the features must be recovered accurately in order to correctly recover camera motion.

Feature correspondence, whether optical flow or discrete features, is based in some form or another on the constant brightness assumption. In its strictest form, the constant brightness constraint assumes that the brightness of the corresponding point does not change between views. Alternatively, we might look at some function of the brightness, possibly a nonlinear function such as brightness edges. The constant brightness constraint is a good approximation for many surfaces in the real world, especially if the motion is small. However, the constant brightness constraint is not strong enough to give us true correspondences. Local measurements, for example, cannot give "optical flow" but only normal flow, the image flow estimates in the direction of the image brightness gradient. This is a particular problem in scenes with long, nearly straight edges. We will show some examples later in the paper.

In this paper, we present the *model-based brightness* constraints where we combine geometric motion models

Recommended for acceptance by P. Flynn.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number 109255.

with the brightness constraint. This provides us with a stronger constraint that can be used for direct estimation of structure and motion. We can threreby bypass the correspondence problem and recover the camera motion directly from the image brightness values. This results in accurate motion estimates and good 3D reconstruction in very challenging scenes.

We now provide the key ideas behind the new method. The *optical flow constraint equation* [17] provides a matching constraint between a point in one image and a line passing through the corresponding point in the second image. In other words, for every point in one image, it gives us the equation of a line along which the corresponding point must lie in the second image. By *point*, we refer to every pixel in the image which has a none zero brightness derivative. The equation is given in terms of the spatiotemporal derivatives of the image brightness and the line is parallel to the iso-brightness contour in the first image. The "optical flow constraint equation" is a first order approximation and assumes that the image motion (and, hence, typically, the camera motion) is small.

There are no geometric constraints on point-line correspondences between two views. This can be shown using the following reasoning: A point in the image defines a line in space. A line in the second image defines a plane in space. A line and a plane always intersect. Therefore, given a set of points in one image and a set of corresponding lines in the second image, for every camera geometry there exists a valid 3D interpretation and no constraint on the geometry exists. We must therefore use three view geometry, where a point in one view and lines through the corresponding points in two other views, provide a constraint, which can be written in the form of the trilinear tensor [31], [13], [39]. The 27 coefficients of the "trilinear tensor" encapsulate the camera motions and the internal parameters of the camera such as the focal length. The new method can therefore be viewed as a necessary extension of the "direct methods" of Horn and Weldon [18] from two views (one motion) to three views (two motions). These methods are dubbed

G.P. Stein is with MobilEye Vision Technology Ltd., Ramot Arazim, 24 Mishof Hadkalim St., Jerusalem, 97278 Israel. E-mail: gideon@moibleye.com.

A. Shashua is with the Hebrew University of Jerusalem, Jerusalem 91904, Israel. E-mail: shashua@cs.huji.ac.

Manuscript received 19 Feb. 1999; revised 11 Aug. 1999; accepted 12 Jan. 2000.



Fig. 1. The three input images (a), (b), and (c) and the estimated depth map (d). The motion between images (a) and (b) is horizontal. The motion is small, but can be seen in the width of the rightmost vertical stripe and in the parallax between the cylinder on the right and the vertical stripes behind it. The motion of between images (a) and (c) was vertical.

"direct methods" because they do not require prior computation of optical flow.

By combining the optical flow constraint equation [17] with the geometric model of the "trilinear tensor" [31], [13], [39], we obtain the tensor brightness constraint [43] that describes the relationship between the spatio-temporal brightness derivatives at each pixel in the image with the camera parameters modeled by the 27 coefficients of the trilinear tensor. This "tensor brightness constraint" provides one linear equation per pixel in the image which results in a highly overconstrained set of equations. The "tensor brightness constraint" is valid for the most general case (projective) where the cameras undergo general motion and we do not know the internal camera parameters which might vary from frame to frame. We then proceed through a hierarchy of reduced motion models, first by assuming calibrated cameras and, then, by assuming the Longuett-Higgins and Prazdny small motion model [21] resulting in reduced model-based brightness constraints for those motion models.

#### 1.1 Related Work

#### 1.1.1 Correspondence-Based Methods

The standard approach to the problem of structure from motion is to first compute correspondences. These might be either dense correspondence in the form of optical flow [17], [22], [26] or feature correspondence [5], [49], [8]. Then, the correspondences are used to compute the camera motion and scene structure. The advantage of our method over both optical flow methods [20], [22], [14], [29], [46] and feature-based methods [47], [8], [37], [38], [1] is that no prior computation of correspondences is needed, a computation which in itself is error prone. Since we obtain a linear system of equations that combines together the information from all the pixels in the image, we avoid the aperture problem without having to apply a smoothness assumption.



Fig. 2. Detail of optical flow computed for Fig. 1. (a) Correct optical flow computed using motion and depth map recovered by the direct method. The camera motion was horizontal parallel to the image plane. The flow vectors have a zero Y component. (b) Optical flow computed using Bergen and Hingorani optical flow program. Note in the upper left corner the flow vectors have a strong Y component. This is due to the aperture problem where even a large window (aperture) will see image gradients in only one direction giving no constraint on the Y component of the flow. Near the intersections (in the image) of the diagonal line and the vertical bars the flow has a small Y component. This error occurs because the program "tracks" the intersection point as if it were a real "feature" point, but the lines do not intersect in space. Although smaller in magnitude, this is the more significant error because the program will also give a high confidence to this value.

By avoiding the need to explicitly detect feature points, we can use information from areas where gradient information is weak such as shadow edges.

These advantages are highlighted in a scene, such as in Fig. 1. The scene contains a background of straight bars, a plaster bust in the foreground, and an oblique line cutting the bars. The background bars provide an unreliable source of information for point-to-point correspondence because of the aperture problem. Even large regions of the image have no information which constrains the flow in the Y direction. Many feature points in the image do not correspond to real feature points in the 3D world. For example, the intersection of the bars and the line in the image do not correspond to physical features in 3D since the line does not lie on the background plane. Fig. 2 shows a detail of the optical flow field in the upper left of the scene where these two problems occur. The flow was computed using a widely available, "industrial strength," optical flow program [3], [2]. (For more examples, see Section 5.3.)

Finally, the texture on the plaster bust varies smoothly and is low contrast thereby giving rise only to a relatively small number of reliable features to track. Yet, our method produces a full 3D model of the scene as demonstrated in Fig. 3a and Fig. 3b. Many natural scenes, such as tree branches or man made objects such as window frames, lamp posts and fences, often give rise to these problems.

#### 1.1.2 Direct Methods

The "direct methods" were pioneered by Horn and Weldon in [18]. Using only a single image pair, they ended up with N equations in N + 5 unknowns, where N is the number of points in the image. The unknowns are the N unknown depths and the three translation and three rotation parameters with one unknown dropping out because translation and depth can only be found up to a common scale factor. The problem is therefore ill-posed and additional constraints are needed. Negahdaripour and Horn [27] present a



Fig. 3. Three-dimensional rendering of the estimated depth map from images in Fig. 1. These images show the inverse depth,  $k = \frac{1}{z}$ , which is the natural value to compute. In (b), the texture map was removed to show the detail and the flaws.

closed form solution assuming a planar or quadratic surface. Szeliski and Kang [45] describe an iterative solution using splines to enforce a smoothness constraint on the depth. McQuirk [23] shows that in a pure translation model the subset of the image points with a nonzero spatial derivative but a zero time derivative gives the direction of motion, thus, the *focus of expansion* (FOE) is on a line perpendicular to the gradient at these points. A significant drawback of this method is that it uses only a small subset of the image points and ignores most of the image data.

Heel [15] constructs a Kalman filter to build up a structure model from more than one image pair, but the core computation is fundamentally the same single image pair computation. The basic idea is that the Horn and Weldon equation is linear in depth if motion is known and linear in motion if depth is known. Heel initially assumes a depth (constant depth surface) and estimates the motion and then uses the motion to estimate the depth. The depth map is then warped using the computed motion and used as the initial depth estimate for the next image pair in the sequence. The question of convergence is not answered and failures are not reported. Results are only shown for pure translation. Heel limits the motion to image motions of less than three pixels which is one of the limiting factors in the accuracy of the surface reconstructions. Within a coarse-to-fine implementation (Section 4.4), our method can handle much larger image motions averaging up to 50 pixels for  $640 \times 480$  resolution images, thereby increasing the dynamic range by an order of magnitude.

Michaels [24] uses three frames (two motions). Each motion gives N equations totaling 2N equations. The unknowns are the N unknown depths and the  $2 \times 6$  translation and rotation parameters. Again, one parameter drops out because of the translation and depth scale ambiguity. He solves the 2N equations with N + 11 unknowns as a large nonlinear optimization problem using the Levenberg-Marquart algorithm [25]. Since it took a very long time to converge, results are shown for very low resolution images only. He also shows theoretically that a large field of view ( $\approx 120^\circ$ ) is required for accurate direct estimation of motion given structure. Since estimation of motion given structure is a key stage in the process, he suggests that a wide field of view is required for the whole process of estimation of structure and motion. Kumar and Anandan [20] first align a dominant plane in the images. The residual image motion is epipolar motion (i.e., pure translation). The epipole of the residual motion is found using iterative techniques starting with an initial guess. Others who have explored planar alignment include [19], [30]. The theoretical framework for planar alignment can also be found in [33].

More recently, Fermuller and Aloimonos [10], [11] describe global geometric properties of the flow-field that give rise to direct relationships between the measurement of normal flow and the ego-motion parameters solved by means of search techniques. They apply the *depth is positive* constraint and try to find an optimally smooth surface.

The tensor brightness constraint was first presented in [36] and a practical implementation with results was first described in [43]. In this paper, we present new theoretical results which lead to a modification of the algorithm. The modified algorithm gives accurate motion estimates even in the presence of considerable camera rotation. We present new quantitative and qualitative results including some with hand-held cameras and outdoor scenes.

# 2 MATHEMATICAL BACKGROUND

## 2.1 Notation

We will use uppercase bold to denote matrices (e.g., **A**). Vectors representing 3D points which will be upper case and not bold. Other vectors and scalars will be in lower case.

We will occasionally use tensorial notations. We use the covariant-contravariant summation convention: A point is an object whose n coordinates are specified with superscripts, i.e.,  $p^i = (p^1, p^2, \ldots, p^n)$ . These are called contravariant vectors. An element in the dual space (representing hyperplanes—lines in  $\mathcal{P}^2$ ), is called a covariant vector and is represented by subscripts, i.e.,  $s_j = (s_1, s_2, \ldots, s_n)$ . Indices repeated in covariant and contravariant forms are summed over, i.e.,  $p^i s_i = p^1 s_1 + p^2 s_2 + \ldots + p^n s_n$ . This is known as a contraction. An outer-product of two one-valence tensors (vectors),  $a_i b^j$ , is a two-valence tensor (matrix)  $c_i^j$  whose i, j entries are  $a_i b^j$ —note that in matrix form  $C = ba^{\top}$ .

Matching image points across three views will be denoted by p, p', p''; the homogeneous coordinates will be referred to as  $p^i, p'^j, p''^k$ , or alternatively as nonhomogeneous image coordinates (x, y), (x', y'), (x'', y'')—hence,  $p^i = (x, y, 1)$ , etc.

We will now consider three perspective views  $\psi$ ,  $\psi'$ , and  $\psi''$  of a 3D scene. Fig. 4 shows a 3D point  $P \in \mathcal{P}^3$  and its image in the three views  $p \in \psi$ ,  $p' \in \psi'$ , and  $p'' \in \psi''$ . Without loss of generality, we can align the 3D world coordinate system with the coordinate system of the first camera:

$$p \cong \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \cong [\mathbf{I}; 0]P. \tag{1}$$

Thus,  $P \cong (x, y, 1, \rho)^{\top}$ . After we have set the first camera coordinate system, the other two camera coordinate systems are in general given by:



Fig. 4. The geometry of three views: Image 1, Image 2, and Image 3 are three views obtained from cameras centered at o, o', and o'', respectively. Point P in the scene projects to image points p, p', and p''. s' and s'' are any image lines passing through the points p' and p'', respectively.

$$p' \cong \begin{pmatrix} x'\\ y'\\ 1 \end{pmatrix} \cong \tilde{\mathbf{A}}P \cong [\mathbf{A};t']P \cong \mathbf{A}p + \rho t'$$
 (2)

$$p'' \cong \begin{pmatrix} x'' \\ y'' \\ 1 \end{pmatrix} \cong \tilde{\mathbf{B}}P \cong [\mathbf{B};t'']P \cong \mathbf{B}p + \rho t''.$$
(3)

The matrices **A** and **B** are homography matrices from Image 1 to Image 2 and to Image 3, respectively, *due to the same plane*  $\pi$ . The vectors t' and t'' are the epipoles, the projection center of camera one projected onto the image planes of the second and thrid cameras, respectively.  $\rho$  is the relative affine depth [34]. It is important to note that  $\rho$  is independent of the second view point and, thus, has the same value in (2) and (3). In a calibrated setting, (i.e., if the intrinsic parameters are known), the matrices **A** and **B** are rotations and the vectors t' and t'' are the translations. Finally, in the calibrated setting,  $\rho$  is replaced by  $k = \frac{1}{2}$ .

We will now derive the triliner tensor of [31]. Let s' and s'' be lines through points p' and p'', respectively:

$$s'^{\top}p' = 0$$
  $s''^{\top}p'' = 0.$  (4)

Premultiplying the left and right hand sides of (2) by  $s'^{\top}$  and (3) by  $s''^{\top}$ , we get:

$$s^{\prime \top} \mathbf{A} p + \rho s^{\prime \top} t^{\prime} = 0$$
  
$$s^{\prime \prime \top} \mathbf{B} p + \rho s^{\prime \prime \top} t^{\prime \prime} = 0.$$
 (5)

Eliminating  $\rho$  from the above equations results in the equation:

$$s'^{\top}t's''^{\top}\mathbf{B}p - s''^{\top}t''s'^{\top}\mathbf{A}p = 0.$$
 (6)

This can be written compactly using tensor notation:

$$p^i s^{\prime\prime}{}_k s^\prime_j \mathcal{T}^{jk}_i = 0, \tag{7}$$

where  $\mathcal{T}$  is the tensor representing a bilinear function of the camera matrices:

$$\mathcal{T}_{i}^{jk} = t'^{j}b_{i}^{k} - t''^{k}a_{i}^{j}.$$
(8)

Equation (7) relates a point p in Image 1 and lines s' and s'' passing through the corresponding points p' and p'' in Image 2 and Image 3, respectively. It is important to note that the lines s' and s'' do not have to correspond to any physical line in space or the image. These constraints first became prominent in [31] and the underlying theory has been studied intensively in [37], [13], [35], [9], [48], [16], [32].

# 3 MODEL-Based BRIGHTNESS CONSTRAINTS

#### 3.1 Photometric Constraints

Geometrically, a trilinear matching constraint is produced by contracting the tensor with the point p in Image 1, *any* line coincident with p' in Image 2, and *any* line coincident with p'' in Image 3. In particular, we may use the tangent to the iso-brightness contour at p' and p'', respectively, and thus one can recover in principle the camera matrices across three views in the context of the "aperture" problem, as suggested by [39]. However, there still remains the problem of finding those matching tangents in the first place. This we now solve.

A first order approximation of the constant brightness constraint leads to the optical flow constraint equation [18]:

$$u'I_x + v'I_y + I'_t = 0, (9)$$

where (u', v') are the *optical flow* values at (x, y) between Image 1 and Image 2 (i.e., u' = x' - x and v' = y' - y).  $(I_x, I_y, I'_t)$  are the spatial and temporal derivatives at the coordinates (x, y). In practice,  $I'_t = I_2(x, y) - I_1(x, y)$ .

The optical flow constraint (9) can be rewritten in the form:

$$(I_x, I_y, I'_t)^{+}(u', v', 1) = 0.$$

The line

$$s = (I_x, I_y, -xI_x - yI_y)^{\top}$$

in the projective plane passes through the point  $p = (x, y, 1)^{\top}$  since:

$$(I_x, I_y, -xI_x - yI_y)^{\top}(x, y, 1) = 0.$$

Combining those two equations together:

$$\begin{pmatrix} I_x \\ I_y \\ -xI_x - yI_y \end{pmatrix}^{\top} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}^{+} \begin{pmatrix} I_x \\ I_y \\ I'_t \end{pmatrix}^{\top} \begin{pmatrix} u' \\ v' \\ 1 \end{pmatrix} = \begin{pmatrix} I_x \\ I_y \\ I'_t - xI_x - yI_y \end{pmatrix}^{\top} \begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix} = 0,$$
(10)

where we have used x' = x + u' and y' = y + v'. Therefore, the line:

$$s' = \begin{pmatrix} I_x \\ I_y \\ I'_t - xI_x - yI_y \end{pmatrix}^{\top}$$
(11)

passes through the point  $p' = (x', y', 1)^{\top}$ .

Thus, the photometric constraints provide a matching constraint between a point p and a line s' passing through the corresponding point p' in Image 2, and between a point p and a line s'' passing through the corresponding point p'' in Image 3. The lines:

$$s' = \begin{pmatrix} I_x \\ I_y \\ I'_t - xI_x - yI_y \end{pmatrix}$$
(12)

and

$$s'' = \begin{pmatrix} I_x \\ I_y \\ I''_t - xI_x - yI_y \end{pmatrix}$$
(13)

are lines coincident with p' and p'', respectively, and parallel to the iso-brightness contour at (x, y).  $I''_t$  is the temporal derivative between the Image 3 and Image 1. (i.e.,  $I''_t = I_3(x, y) - I_1(x, y)$ .)

## 3.2 The Projective Model: The Tensor Brightness Constraint

Substituting (12) and (13) into the tensor (8) results in the *tensor brightness constraint*:

$$s''_{k}s'_{j}p^{i}\mathcal{T}_{i}^{jk} = 0.$$
(14)

We have one such equation for each point on the image where  $s''_k$  and  $s'_j$  can be computed from the image gradients and  $p^i = (x, y, 1)$  are the (projective) image coordinates of the point in Image 0. We solve for  $\mathcal{T}_i^{jk}$ which combines the motion and camera parameters. The coordinates of the corresponding points (x', y') and (x'', y'') are not required.

Every pixel with a nonvanishing gradient contributes one linear equation to the 27 unknown parameters comprising  $\mathcal{T}_i^{jk}$ . However, a configuration of a point in the first image and two *parallel* lines in Images 2 and 3 is a particular instance of a degenerate line configuration called a Linear Line Complex. The general solution for the LLC case is explored in [44]. In this particular case, the system of equations can provide a linear solution to 21 of the parameters and the remaining six parameters of  $\mathcal{T}_i^{jk}$  can be determined using quadratic admissibility constraints. Here, we provide a proof for this special case:

From (14), the coefficients of the terms  $\mathcal{T}_i^{12}$  and the terms  $\mathcal{T}_i^{21}$  are both  $p^i I_x I_y$ . Therefore, the linear equations cannot be solved for the six terms ( $\mathcal{T}_i^{12}$  and  $\mathcal{T}_i^{21}$ ) individually, but only for the three sums

$$\mathcal{T}_{i}^{12} + \mathcal{T}_{i}^{21}, i = 1 \dots 3.$$

The sum  $\mathcal{T}_i^{12} + \mathcal{T}_i^{21}$  provides one linear equation in  $\mathcal{T}_i^{12}$  and  $\mathcal{T}_i^{21}$ . We will use the admissibility constraint on the Standard Correlation Slices,  $T_i$ . The constraint states that the matrix  $T_i$  is of rank = 2. This leads to a quadratic equation in  $\mathcal{T}_i^{12}$  and  $\mathcal{T}_i^{21}$ . The quadratic and linear equations together result in two solutions for each pair  $\mathcal{T}_i^{12}, \mathcal{T}_i^{21}, i = 1...3$  for a total of eight discrete solutions. A unique solution is obtainable by applying further admissibility constraints, as shown in [44].

This added complexity in finding the unique solution leads us to search for simpler models. Next, we consider the case of a calibrated camera with the rotation limited to small rotation angles.

## 3.3 The Small Rotation Model with Calibrated Cameras

The next model is defined for small-angle rotations with calibrated cameras. Assuming calibrated cameras (2) and (3) become:

$$p' \cong [\mathbf{R}'; t']P \cong \mathbf{R}'p + kt'$$
  
$$p'' \cong [\mathbf{R}''; t'']P \cong \mathbf{R}''p + kt''$$
(15)

i.e.,  $\tilde{A} = [\mathbf{R}'; t']$  and  $\tilde{B} = [\mathbf{R}''; t'']$  and  $k = \frac{1}{z}$  replaces  $\rho$  in (2) and (3).

If we also assume small angle rotations and can therefore make the approximations:

$$\cos(\theta) \approx 1, \ \sin(\theta) \approx \theta$$

then the rotation matrices can be approximated as:

$$R' \approx \left[ I + [w']_x \right] R'' \approx \left[ I + [w'']_x \right],$$
(16)

where w', w'' are the angular velocity vectors and  $[\cdot]_x$  is the skew-symmetric matrix of vector products. Now, (2) and (3) become:

$$p' \cong \tilde{\mathbf{A}}P \cong [I + [w']_x; t']P$$

$$p'' \cong \tilde{\mathbf{B}}P \cong [I + [w'']_x; t'']P.$$
(17)

We can now proceed along two paths. The different results will give us some further insight into the problem. First, we will substitute  $\tilde{A} = [I + [w']_x; t']$  and  $\tilde{B} = [I + [w']_x; t'']$  directly into (8) and then write out the tensor explicitly:

$$T_{1} = \begin{bmatrix} t_{1}' - t_{1}'' & t_{1}'w_{3}'' - t_{2}'' & -t_{1}'w_{2}'' - t_{3}'' \\ t_{2}' - t_{1}''w_{3}' & t_{2}'w_{3}'' - t_{2}'w_{3}' & -t_{2}'w_{2}'' - t_{3}'w_{3}' \\ t_{3}' + t_{1}''w_{2}' & t_{3}'w_{3}' + t_{2}''w_{2}' & -t_{3}'w_{2}' + t_{3}''w_{2}'' \end{bmatrix}$$
(18)

$$T_{2} = \begin{bmatrix} -t'_{1}w''_{3} + t''_{1}w'_{3} & t'_{1} + t''_{2}w'_{3} & t'_{1}w''_{1} + t''_{3}w'_{3} \\ -t'_{2}w''_{3} - t''_{1} & t'_{2} - t''_{2} & t'_{2}w''_{1} - t''_{3} \\ -t'_{3}w''_{3} - t''_{1}w'_{1} & t'_{3} - t''_{2}w'_{1} & t'_{3}w''_{1} - t''_{3}w'_{1} \end{bmatrix}$$
(19)

$$T_{3} = \begin{bmatrix} t_{1}'w''_{2} - t''_{1}w'_{2} & -t_{1}'w''_{1} - t''_{2}w'_{2} & t_{1}' - t''_{3}w'_{2} \\ t_{2}'w''_{2} + t''_{1}w'_{1} & -t_{2}'w''_{1} + t''_{2}w'_{1} & t_{2}' + t''_{3}w'_{1} \\ t_{3}'w''_{2} - t''_{1} & -t_{3}'w''_{1} - t''_{2} & t_{3}' - t''_{3} \end{bmatrix}, \quad (20)$$

where  $T_1, T_2$ , and  $T_3$  are the Standard Correlation Slices [46]. Studying the tensor terms, we notice that the terms are not linearly independent and the following seven linear equations hold:

$$\begin{split} \mathcal{T}_{1}^{1,2} + \mathcal{T}_{1}^{1,1} + \mathcal{T}_{1}^{2,1} - \mathcal{T}_{2}^{2,2} &= 0 \\ \mathcal{T}_{2}^{2,1} + \mathcal{T}_{1}^{2,2} + \mathcal{T}_{2}^{1,2} - \mathcal{T}_{1}^{1,1} &= 0 \\ \mathcal{T}_{3}^{1,1} + \mathcal{T}_{1}^{1,3} + \mathcal{T}_{3}^{3,1} - \mathcal{T}_{3}^{3,3} &= 0 \\ \mathcal{T}_{1}^{3,3} + \mathcal{T}_{3}^{3,1} + \mathcal{T}_{3}^{3,2} - \mathcal{T}_{1}^{1,1} &= 0 \\ \mathcal{T}_{2}^{3,3} + \mathcal{T}_{2}^{3,2} + \mathcal{T}_{2}^{2,3} - \mathcal{T}_{2}^{2,2} &= 0 \\ \mathcal{T}_{3}^{2,2} + \mathcal{T}_{2}^{2,3} + \mathcal{T}_{2}^{3,2} - \mathcal{T}_{3}^{3,3} &= 0 \\ \mathcal{T}_{2}^{1,3} + \mathcal{T}_{3}^{1,1} + \mathcal{T}_{1}^{2,3} + \mathcal{T}_{3}^{1,2} + \mathcal{T}_{1}^{3,2} + \mathcal{T}_{2}^{3,1} &= 0. \end{split}$$

$$(21)$$

Therefore, for the calibrated small rotation model, the 27-parameter Tensor Brightness Constraint reduces to a

20-parameter constraint equation. Each point in the image gives one homogeneous equation resulting in a set of N homogeneous equations, where N is the number of image points. This can be written in matrix form as:

$$\mathbf{A}x = 0$$

where **A** is the  $N \times 20$  estimation matrix and x is the 20-parameter vector of the unknown intermediate parameters. There is a unique nontrivial solution if the estimation matrix is of rank = 19. But, as in the projective case, the system of homogeneous equations is degenerate and, in this case, the rank of the estimation matrix is 16 not 19.

In [36], [43], a model-based brightness constraint for the calibrated, small rotation model was derived in a different way. Expanding (17):

$$p' \cong \left[I + [w']_x\right]p + kt' \tag{22}$$

and taking the dot product of the left-hand-side and righthand-side terms with s' (12) results in the equation:

$$0 = s'^{\top} p + s'^{\top} [w']_x p + k s'^{\top} t'.$$
(23)

This can be simplified to the form:

$$ks'^{\top}t' + q'^{\top}w' + I'_t = 0 \tag{24}$$

by noting that  $s'^{\top}p = I'_t$  and  $s'^{\top}[w']_x p = q'^{\top}w'$  if we define:

$$q' = p \times s' = \begin{pmatrix} -I_y + y(I'_t - xI_x - yI_y) \\ I_x - x(I'_t - xI_x - yI_y) \\ xI_y - yI_x \end{pmatrix}.$$
 (25)

In a similar manner, for the second motion:

$$ks''^{\top}t'' + q''^{\top}w'' + I''_{t} = 0, \qquad (26)$$

where

$$q'' = p \times s'' = \begin{pmatrix} -I_y + y(I''_t - xI_x - yI_y) \\ I_x - x(I''_t - xI_x - yI_y) \\ xI_y - yI_x \end{pmatrix}.$$

Multiplying (24) by  $s''^{\top}t''$  and (25) by  $s'^{\top}t'$  and subtracting, we obtain the following 24-parameter model-based brightness constraint as a reduction of the tensor brightness constraint:

$$I''_{t}s'^{\top}t' - I'_{t}s''^{\top}t'' + s'^{\top}[t'w''^{\top}]q'' - s''^{\top}[t''w'^{\top}]q' = 0.$$
 (26)

The resulting 24-parameter model (26) is not as compact as the 20-parameter model we obtained by direct substitution of the small rotation approximation into the tensor brightness constraint equation. We find that the  $24 \times N$  estimation matrix obtained using the 24-parameter model is again of rank = 16. Due to the high degeneracy of the linear equations, neither the 24-parameter model nor the 20-parameter model are convenient to work with.

#### 3.4 The 15-Parameter, Small Motion Model

A unique solution is obtained when we reduce the motion model further to include infinitesimal motion using the model introduced by Longuet-Higgins and Prazdny [21]. The LH&P model assumes in addition to small rotation that  $\frac{t_z'}{Z} \ll 1$ . The motion field equations for the first camera motion are then:

$$u' = \frac{1}{z}(t'_1 - xt'_3) - w'_3 y + w'_2(1 + x^2) - w'_1 x y$$
  

$$v' = \frac{1}{z}(t'_2 - yt'_3) + w'_3 x - w'_1(1 + y^2) + w'_2 x y.$$
(27)

By substituting (27) into the optical flow constraint equation:

$$u'I_x + v'I_y + I'_t = 0 (28)$$

and rearranging the terms, we obtain:

$$ks^{\top}t' + q^{\top}w' + I'_t = 0, \qquad (29)$$

where  $k = \frac{1}{z}$  denotes the inverse depth at each pixel location and where s, q are defined below:

$$s = \begin{pmatrix} I_x \\ I_y \\ -xI_x - yI_y \end{pmatrix}$$
(30)

and

$$q = p \times s = \begin{pmatrix} -I_y - y(xI_x + yI_y) \\ I_x + x(xI_x + yI_y) \\ xI_y - yI_x \end{pmatrix}.$$
 (31)

Equation (29), which was first derived in [27], [18], can also be derived directly from (24) by making the LH&P assumptions:  $\frac{u'_x}{Z} \ll 1$ ,  $\frac{w'_y x}{f} \ll 1$ , and  $\frac{w'_x y}{f} \ll 1$ .

Similarly, for the second motion:

$$ks't'' + q'w'' + I''_t = 0 (32)$$

By eliminating *k* from (29) and (32), we obtain the 15-parameter model-based brightness constraint:

$$I''_{t}s^{T}t' - I'_{t}s^{T}t'' + s^{T}[t'w''^{T} - t''w'^{T}]q = 0.$$
 (33)

Equation (33) is a key equation. There is one such equation for every point in the image. The unknowns are the motion parameters t', t'', w', and w''. The values  $I'_t$ ,  $I''_t$ , s, and v, can all be computed from image derivatives and the coordinates of the point in Image 1. Given the solution to the ego-motion parameters (to be described later), one can recover the dense depth map from (29) and (32).

Before showing how to solve (33) for both rotation and translation, we will start with the simple case of pure translation. This case arises, in practice, when we can either assure pure translation or we have previously rectified the images (perhaps by registering a common plane in the images as in [20], [19], [30]). It is also an important case from a theoretical point of view. We will show that there are problems in recovering the motion parameters in the case of collinear motion. If collinear motion causes a problem for the pure translation case, it will also create a problem for the more general translation and rotation case.

#### 3.5 Pure Translation

The pure translational small-motion model takes the simplest form:

998

$$I''_{t}s^{T}t' - I'_{t}s^{T}t'' = 0. (34)$$

We have one such equation for each image point and we can write it out in the matrix form:

$$\mathbf{A}t = 0,$$

where:

$$t = \begin{pmatrix} t'_{x} & t'_{y} & t'_{z} & t''_{x'} & t''_{y'} & t''_{z'} \end{pmatrix}^{T}.$$

and **A** is an  $N \times 6$  matrix with the *i*th row (corresponding to the *i*th pixel) given by:

$$(I''_t s_{i1} \quad I''_t s_{i2} \quad I''_t s_{i3} \quad -I'_t s_{i1} \quad -I'_t s_{i2} \quad -I'_t s_{i3}).$$

To avoid the trivial solution t = 0, we add the constraint ||t|| = 1. The least-squares problem now maps to the problem of finding ||t|| = 1 that minimizes:

$$t^T \mathbf{A}^T \mathbf{A} t = 0.$$

The solution is the eigenvector of  $\mathbf{A}^T \mathbf{A}$  corresponding to the smallest eigenvalue.

#### 3.5.1 The Singularity of Collinear Motion

This method fails when the two motions are in the same (or opposite) directions. In the pure translation case,

$$ks't' + I'_t = 0 (35)$$

and

$$ks^{\top}t'' + I''_{t} = 0. \tag{36}$$

If the second translation vector t'' is proportional to the translation vector t' then (36) is simply a scaled version of (35) adding no new information. The solution will, therefore, be ill-conditioned when  $t' \approx \beta t''$  for a scalar  $\beta$ . This is a drawback in many applications (e.g., 3D reconstruction from a monocular image sequence). In Section 6.1, we present some early research on ways to overcome this problem. The initial results look promising.

#### 3.6 Solving for Translation and Rotation

In (33), the ego-motion parameters are embedded in a 15-parameter model: six translation parameters and the nine outer-product terms  $[t'w''^T - t''w'^T]$ . Each pixel provides one linear equation, thus N pixels provide a system  $\mathbf{A}c = 0$ , where  $\mathbf{A}$  is an  $N \times 15$  measurement matrix and c is the 15-vector of unknown parameters. The least-squares solution is the eigenvector of  $\mathbf{A}^{\top}\mathbf{A}$  corresponding to the smallest eigenvalue. Without noise, the eigenvalue will be zero and the solution is in the null space of  $\mathbf{A}^{\top}\mathbf{A}$ . The solution is unique if the null space is one-dimensional and the rank of A is 14. We will prove, however, that the null space is at least two-dimensional.

**Proposition 1.** Let **A** be the  $N \times 15$  measurement matrix associated with the homogeneous (33), where  $N \gg 15$ . Then the null space of the matrix  $\mathbf{A}^{\top}\mathbf{A}$  is of dimension greater or equal to 2 and  $Rank(\mathbf{A}) \leq 13$ .

Proof. Since

$$s^\top q = s^\top (p \times s) = 0$$

the vector

$$c_0 = (0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1)^3$$

is in the null space of **A**.

The vector  $c_0$  corresponds to:

$$\begin{aligned} t' &= 0\\ t'' &= 0\\ \begin{bmatrix} t'w''^T - t''w'^T \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0\\ 0 & 1 & 0\\ 0 & 0 & 1 \end{bmatrix}. \end{aligned}$$

Therefore,  $c_0$  is clearly not an admissible solution and, therefore, the null space of **A** includes two vectors: the vector  $c_0$  and the true solution. It follows that the rank of **A** is at most 13.

The method of solution is based on the following observation. Let  $b_0$  be another vector of the null-space (for example, let  $c_0$  and  $b_0$  be the two eigenvectors of  $\mathbf{A}^{\top}\mathbf{A}$  corresponding to the two smallest eigenvalues), then the desired solution vector b is a linear combination of the two:

$$b = b_0 - \alpha c_0. \tag{37}$$

We note that  $c_0$  is in the null space of **A** regardless of noise in the data and will, therefore, correspond to an eigenvalue of zero, up to numerical precision. Theoretically,  $b_0$  is also in the null space of **A**, but only in the case of noise free data and exact motion model. Therefore,  $b_0$  is the eigenvector corresponding to the second smallest eigenvalue. Next, we show how to determine  $\alpha$ .

## 3.6.1 Solution Using the Rank = 2 Constraint

In order to find  $\alpha$  given  $c_0$  and  $b_0$ , we enforce the constraint that:

$$Rank[t'w''^{T} - t''w'^{T}] = 2.$$

Clearly, the choice of  $\alpha$  will have no affect in the first six elements of the vector *b*. Let us arrange the last nine elements of *b*, *b*<sub>0</sub>, and *c*<sub>0</sub> into the corresponding  $3 \times 3$  matrices **B**, **B**<sub>0</sub>, and **C**<sub>0</sub>. We are now looking for an  $\alpha$  such that:

$$Rank(\mathbf{B}_0 - \alpha \mathbf{C}_0) = 2. \tag{38}$$

In our case,  $C_0$  is the identity matrix, so (38) becomes:

$$Rank(\mathbf{B_0} - \alpha \mathbf{I}) = 2$$

and the solution for  $\alpha$  is given simply by the eigenvalues of **B**<sub>0</sub>. Since **B**<sub>0</sub> is a 3 × 3 matrix, this results in up to three discrete solutions. We now prove that only one of those can be a valid solution.

**Theorem 1.** The 15-parameter model has a unique solution. Let:

$$\mathbf{B} = \begin{bmatrix} t'w''^T - t''w'^T \end{bmatrix}.$$
(39)

Let:

$$\tilde{\mathbf{B}} = \left[t'w''^T - t''w'^T\right] + \alpha \mathbf{I}.$$
(40)

*Then, the equation:* 

$$\tilde{\mathbf{B}} = \begin{bmatrix} t'\tilde{w}''^T - t''\tilde{w}'^T \end{bmatrix}$$
(41)

has a solution for  $\tilde{w}'', \tilde{w}'$  given t', t'' iff  $\alpha = 0$ .

**Proof.** The *if* part of the proof is trivial. For  $\alpha = 0$ ,  $\tilde{w}'' = w''$ , and  $\tilde{w}' = w'$  is a solution for (41). To prove the *only if* part, suppose there exist  $\tilde{w}'', \tilde{w}'$  that is a solution for (41). Then,

$$[t'w''^{T} - t''w'^{T}] + \alpha \mathbf{I} = [t'\tilde{w}''^{T} - t''\tilde{w}'^{T}].$$
(42)

Let *q* be a vector such that  $q^T t' = 0$ ,  $q^T t'' = 0$ , and  $q \neq \vec{0}$ . Multiply the left and right side of (42) by  $q^T$  to get:

$$q^{T}[t'w''^{T}] - q^{T}[t''w'^{T}] + q^{T}\alpha \mathbf{I} = q^{T}[t'\tilde{w}''^{T}] - q^{T}[t''\tilde{w}'^{T}].$$
(43)

Since  $q^T t' = 0$  and  $q^T t'' = 0$ , most of the terms in (43) become zero and we are left with:

$$q^T \alpha \mathbf{I} = \vec{0}. \tag{44}$$

Since  $q \neq \vec{0}$ , then  $\alpha = 0$ .

#### 3.6.2 The Algorithm for Finding Motion Parameters

Based on the previous arguments, the algorithm for finding the motion parameters is as follows:

- 1. Compute the  $N \times 15$  matrix A from (33).
- 2. Find the eigenvector corresponding to the second smallest eigenvalue of the matrix  $\mathbf{A}^T \mathbf{A}$ . This is  $b_0$ . (The vector  $c_0$  corresponds to the smallest.)
- 3. The first six elements of  $b_0$  are the translations, t', t''.
- 4. Arrange the last nine elements of  $b_0$  into a  $3 \times 3$  matrix **B**<sub>0</sub>.
- 5. Construct three possible rank two matrices  $\mathbf{B}_{i,i=1..3}$  from  $\mathbf{B}_0$ :

$$\mathbf{B}_i = \mathbf{B}_0 - \alpha_i \mathbf{I},\tag{45}$$

where  $\alpha_{i,i=1..3}$  are the three eigenvalues of **B**<sub>0</sub>. For each  $\alpha_{i,i=1..3}$  solve:

6.

$$[t'w''^T - t''w'^T] = \mathbf{B},$$

for w' and w'', given t' and t'' from Step 3. This a set of nine linear equations in the six unknowns (w', w'') and is solved using least squares.

From Theorem 1, only one of the three solutions is correct. We select the solution that best fits the data. For each of the three solutions, use the 12-parameters t', t", w', and w'' to form the 15 intermediate parameter vector b<sub>i</sub> and then compute the error:

$$E_i = b_i^T \mathbf{A}^T \mathbf{A} b_i. \tag{46}$$

Of the three solutions, select the solution which gives the smallest error.

#### 3.6.3 Solve as a Nonlinear Optimization

In Section 3.6.2, we use linear methods to solve for the motion parameters. In order to do so, we have treated the 15 intermediate parameters as linearly independent parameters while they are in fact bilinear combinations of 12 independent parameters. It is often possible to improve the motion estimates by using the linear solution

as a starting point for nonlinear optimization techniques [7]. In our experiments, we did not find that a nonlinear optimization stage improved the results. This is possibly because the iterative framework (Section 4.2) is in itself a form of nonlinear optimization.

The nonlinear optimization procedure is brought here for completeness. Our general problem is to find motion parameters t', t'', w', and w'' that minimize the cost function:

$$E(t', t'', w', w'') = b^{\top} \mathbf{A}^{\top} \mathbf{A} b, \qquad (47)$$

where **A** is the estimation matrix derived from (33) and b(t', t'', w', w'') is the fifteen element vector of bilinear functions of the motion parameters.

One method would be to find motion parameters t', t'', w', and w'' that minimize (47). In order to avoid the trivial solution, we need to add a constraint such as:  $|t'|^2 + |t''|^2 = 1$ . This becomes a 12-parameter nonlinear constrained optimization problem.

The process can be simplified by noting that (33) can be rewritten in the form:

$$(I''_t - w''^{\top} q'') s'^{\top} t' - (I'_t - w'^{\top} q') s''^{\top} t'' = 0$$
(48)

given a set of rotation values w' and w'', one can compute the least-squares estimate of t' and t'' using linear methods and also compute the least-squares error (47). One can then use nonlinear search techniques to find w' and w'' which minimize this least-squares error. This becomes a six parameter unconstrained optimization problem.

#### 4 IMPLEMENTATION DETAILS

# 4.1 Computing the Depth, Smoothing, and Interpolation

After recovering the camera motion (Section 3.6.2), we use (29) and (32) to compute the depth at every point. Information is combined from both image pairs by minimizing the least-squares error:

$$E = \min_{k} \arg \sum_{j=1}^{2} \left( k s^T t^j + q^T w^j + I_t^j \right)^2.$$

Here, j = 1 and j = 2 denote values from the first and second image pairs, respectively.

There are points in the image where the brightness gradients are close to zero (and, therefore,  $s^T t^j \simeq 0$ ) and the estimation of k will be ill-conditioned at those points. In order to overcome this problem, we use a local region of support around the point of interest. In it's simplest form, we assume the depth is constant in the region and minimize:

$$E = \min_{k} \arg \sum_{x,y \in R} \sum_{j} \beta(x,y) |s^{T}t^{j}|^{p} \left( ks^{T}t^{j} + q^{T}w^{j} + I_{t}^{j} \right)^{2},$$
(49)

where the windowing function  $\beta(x, y)$  allows one to increase the weights of the points closer to the center of the region. The  $|s^T t^j|^p$  term reduces the weight of points which have a small gradient or where the gradient is

perpendicular to the camera motion since these cases are highly affected by noise. We used p = 1.

During the iteration process, we typically used a region R of  $7 \times 7$  to  $11 \times 11$ . For "prettier" results in the last iteration, we typically reduced the region R to  $1 \times 1$  but added a very weak global smoothness term and performed multigrid membrane interpolation [28]. The smoothness term stabilizes regions where there is no image gradient so very small regions of support can be used.

#### 4.2 Iterative Refinement

The constant brightness constraint is a linearized form of the Sum Square Difference (SSD) criterion. The linear solution can be thought of as a single iteration of Newton's method applied to the problem. Iterative refinement is performed as follows:

- 1. Calculate motion (using (33)).
- 2. Compute depth (using (49)).
- 3. Using the depth and motion, warp Images 2, and 3 towards Image 1 (Section 4.3).
- 4. Compute new time derivatives  $I'_t$  and  $I''_t$ .
- 5. Compute a new motion and depth estimate.

In the ideal case, as the final result, the warped images should appear nearly identical to Image 1. One must be careful and not simply compute the incremental changes in the translation  $\delta t'$  and  $\delta t''$ . As the images are warped closer together and the translation estimate approaches zero, the system of (33) will become ill-conditioned. Furthermore, since the equations are homogeneous, we must enforce a constraint such as  $|t'|^2 + |t''|^2 = 1$  to avoid the trivial solution. We would not wish to apply such a constraint to  $\delta t'$  and  $\delta t''$ . Thus, one must compute the full translation model (previous iteration plus the incremental change to the translation). This problem does not arise for the rotations: w' and w''. It is, in fact, convenient to warp the images using the best rotation estimate and then compute only the incremental improvement in the rotations:  $\delta w'$  and  $\delta w''$ . For example, we might choose to warp the images using a more exact rotation model than the linearized small rotation model.

The method for computing the full translation model and incremental rotation is now described. Let  $\Psi_1$ ,  $\Psi_2$ , and  $\Psi_3$  be the three images. Assume we have  $\hat{k}$ ,  $\hat{t}^j$ ,  $\hat{w}^j$ , from the previous iteration. (Here, again, we use *j* to stand for either ' or ".) The translation components of image motion  $\hat{u}_t^j$  and  $\hat{v}_t^j$ , and the rotational components  $\hat{u}_r^j$  and  $\hat{v}_r^j$  can be computed using  $\hat{k}$ ,  $\hat{t}^j$ , and  $\hat{w}^j$  in (27). Then, these are used to warp images  $\Psi_2$  to  $\hat{\Psi}_2$  and  $\Psi_3$  to  $\hat{\Psi}_3$ . After warping, the images satisfy the brightness constraint equation:

$$I_x du' + I_y dv' + I'_t = 0$$
  

$$I_x du'' + I_y dv'' + \hat{I}''_t = 0,$$
(50)

where the temporal derivatives at each pixel are given by:

$$I'_{t} = \Psi_{2} - \Psi_{1}$$

$$\hat{I}''_{t} = \hat{\Psi}_{3} - \Psi_{1}$$
(51)

and  $du^{j}$ ,  $dv^{j}$  are the (still unknown) differences between computed image motions and the real image motions:

$$du^{j} = u^{j} - \hat{u}_{t}^{j} - \hat{u}_{r}^{j} dv^{j} = v^{j} - \hat{v}_{t}^{j} - \hat{v}_{r}^{j}.$$
(52)

Let:

$$\alpha^j = I_x \hat{u}_t^j + I_y \hat{v}_t^j, \tag{53}$$

which can also be written as:

$$\alpha^j = \hat{k} s^T \hat{t}^j. \tag{54}$$

Substituting (53) and (52) in (50), we get:

$$I_x(u' - \hat{u}'_r) + I_y(v' - \hat{v}'_r) + (I'_t - \alpha') = 0$$
  

$$I_x(u'' - \hat{u}''_r) + I_y(v'' - \hat{v}''_r) + (\hat{I}''_t - \alpha'') = 0.$$
(55)

Substituting (27) in (55), we get modified versions of the (29) and (32)

$$ks^{T}t' + q^{T}\delta w' + (I'_{t} - \alpha') = 0$$
  

$$ks^{T}t'' + q^{T}\delta w'' + (I''_{t} - \alpha'') = 0.$$
(56)

We start our first iteration with  $\hat{k}$ ,  $\hat{t}^{j}$ ,  $\hat{w}^{j}$  all zero and, therefore,  $\alpha = 0$  as well.

## 4.3 Image Warping

Given an estimate of the camera motion and the depth at every point, we can warp Image 2 towards Image 1. Image warping, in general, is described in [50]. Let  $I'_{old}$  be the original Image 2 and let  $I'_{new}$  be the warped image we are trying to create. We must first define functions  $x_{old}(x, y, ...)$ and  $y_{old}(x, y, ...)$  which, given image coordinates in the new image (x, y) and possibly some extra parameters (symbolized by ,...), return the coordinates of the corresponding point in the old image  $(x_{old}, y_{old})$ .

In our case, the *extra parameters* are the camera motions t' and w' and the depth map at every point k(x, y). The functions  $x_{old}$  and  $y_{old}$  depend on the motion model. In our implementation, we used the LH&P small motion model (27). Therefore,

$$x_{old} = x + u'$$
  
=  $x + \frac{1}{z}(t'_1 - xt'_3) - w'_3 y + w'_2(1 + x^2) - w'_1 x y$  (57)

and

$$y_{old} = y + v'$$
  
=  $y + \frac{1}{z}(t'_2 - yt'_3) + w'_3 x - w'_1(1 + y^2) + w'_2 xy.$  (58)

After defining  $x_{old}$  and  $y_{old}$ , we can compute the value for every pixel in  $I'_{new}$  according to the formula:

$$I_{new}'(x,y) = \begin{cases} I_{old}'(x_{old}(x,y,\ldots),y_{old}(x,y,\ldots)), \\ if(x_{old},y_{old}) \in I_{old}' \\ 0, & otherwise. \end{cases}$$
(59)

There are two points to note:

• The values of  $x_{old}$  and  $y_{old}$  are, in general, noninteger so we use bilinear interpolation to compute the appropriate pixel values. Bilinear interpolation works better than nearest-neighbor. The use of the more complex bicubic interpolation resulted in no noticeable improvement over bilinear interpolation.

• If the coordinates (*x*<sub>old</sub>, *y*<sub>old</sub>) are outside the coordinates of the image, then, we use the value 0. These points should be marked as invalid and not used in subsequent motion and depth computation.

# 4.4 Coarse-to-Fine Processing

In order to deal with image motions larger than one pixel, we use a Gaussian pyramid for coarse to fine processing [2], [4]. Each of the three images is filtered by an approximation to a Gaussian filter and subsampled to create an image of half the size (in each dimension). This operation is performed recursively to create a pyramid with four or five levels depending on the size of the original image. The following  $5 \times 5$  kernel was used as the filter:

	[1	4	6	4	1	
$\frac{1}{256} \times$	4	16	24	16	4	
	6	24	36	24	6	
	4	16	24	16	4	
	1	4	6	4	1	

This filter is separable and can be implemented by convolving the rows and columns of the image, twice each, with the filter f = (0.25, 0.5, 0.25).

Starting from the coarsest level, we perform a few iterations to compute motion and depth according to the scheme described in Section 4.2. The brightness derivatives  $(I_x, I_y, I'_t, I''_t)$  are computed using the subsampled images. After computing the motion and depth, we create a finer depth map from the coarse depth map using bilinear interpolation. We then use the interpolated depth map and the motion estimates as the starting values for the iterations at the finer level.

#### 4.5 Field of View

In this section, we will first demonstrate how, in the general problem of determining camera motion from motion fields, when the field of view is narrow, there exists an ambiguity between rotation and translation. We will also note the errors in the estimation of the rotation if we use the wrong focal length value. Then, we show how this ambiguity is observed in the shape of the cost function (47) from three views. Finally, we describe how to achieve stable and accurate motion estimates in the presence of rotation even without very wide fields of view.

#### 4.5.1 Rotation-Translation Ambiguity from Motion Fields

It is well-known that for a medium to narrow field of view, the motion field due to camera rotation around the Y axis, can be indistinguishable from the motion field due to some translation along the X axis and an appropriate depth surface. A similar ambiguity exists between rotations around the X axis and translations along the Y axis.

In Fig. 5, we see that for medium-narrow fields of view  $(30^\circ)$ , the flow due to rotation around the *Y* axis is parallel to the *X* axis. The length of the flow vectors due to rotation is slightly shorter, closer to the *Y* axis. A similar flow field can be produced by camera translation along the *X* axis with a surface that is curved so that points in the center of the image are more distant. There are motion flows that



Fig. 5. Motion flow field due to rotation around the *Y* axis for wide, medium, and narrow fields of view. The arrows for the wide field of view (short, black arrows) and the medium field of view (medium length, pale short arrows) have been offset in the Y direction for clarity. For the narrow field of view ( $f = 800_{pixels}$ , FOV =  $30^{\circ}$ ), the flow field is large and almost parallel to the *X* axis. For medium fields of view (f = 400) and wide fields of view (f = 200), the *X* component of the flow is smaller but there is a noticeable *Y* component.

unambiguously indicate a translational component to the motion: motions fields that include a focus-of-expansion or show parallax effects. Nevertheless, even these fields are ambiguous since a rotational component could be added to the motion, compensated for by a change in the translation and depth to produce the same motion field. In some cases, the resulting depth map will include negative depths indicating some obvious error. This is often formulated as the *depth positive constraint* [12] which can be used to limit the range of motion ambiguity.

With wider fields of view (i.e., shorter focal lengths), the ambiguity disappears. The motion fields, due to rotation, cannot be modeled by translational motion. Fig. 6a also shows the motion fields due to rotation with fields of view of  $53^{\circ}$  and  $90^{\circ}$ . The distribution of the *Y* component of the flow has the unique characteristics of rotational flow. For wide fields of view, the *Y* component points towards the *X* axis on the right of the image and away from the X axis on the left of the image (for this particular sign of the rotatation).

The motion field for a wide field of view is unambiguous if we know the focal length, or equivalently, the field of view (i.e., the camera's internal parameters are known). If our estimate of the focal length is incorrect, then the motion and structure estimates which best fit the flow field will be incorrect. If the focal length estimate is smaller than the true focal length, then the estimated angle of rotation will be smaller than the true rotation angle (in absolute values). Fig. 6a shows the Y component of the motion field for a rotation of  $-4.5^{\circ}$  around the Y axis using a focal length of  $400_{pixels}$  and the same measurements for a focal length of  $200_{pixels}$  and a rotation of  $-2.5^{\circ}$ . Fig. 6b shows the X component of the flow. The Y components are similar. The X components are very different, but these can be accounted for by adjusting the translation and depth estimates.



Fig. 6. The *X* and *Y* components of the motion field due a rotation around the *Y* axis of  $-4.5^{\circ}$  with a focal length of  $400_{pixels}$  and for a rotation of  $-2.5^{\circ}$  with a focal length of  $200_{pixels}$ . The arrows for the  $200_{pixels}$  camera have been plotted offset right in (a) and up in (b). Note that the *Y* flow components are very similar.

# 4.5.2 Rotation-Translation Ambiguity in the 15-Parameter Model

The rotation translation ambiguity can also be observed while using the 15-parameter model. Fig. 7 shows the shape of the cost function (47) near the global minimum. These are simulation results for focal lengths of 50, 100, and 200. The image size was  $320 \times 240$ . It is not possible to plot the full 12-dimensional surface so the figures show only the shape of the function as we vary two motion parameters at a time. Figs. 7d, 7e, and 7f show the effect of varying  $w'_x$  and  $t'_{y'}$ , where we expect to see ambiguity and the effect of varying  $w'_x$  and  $t'_x$ , where no ambiguity is expected. In Fig. 7f, the long diagonal valley shape of the cost function indicates the ambiguity of the motion estimate for narrower fields of view.

## 4.5.3 The Solution

As we have seen, in order to get reasonable rotation and translation estimates, a wide field of view is required. The theoretical results of [24] and simulation results indicate that very wide fields of view are required ( $120^{\circ}$  or greater).

These are hard to achieve with standard lenses. With narrower fields of view, the iterative scheme described in Section 4.2 fails to converge. In previously published work [43], we managed to stabilize the results for small rotation angles by unknowingly biasing the results towards small rotation. This was done by setting the focal length parameter to 50 instead of around 600. For a  $640 \times 480$  size image, this implied a field of view of  $150^{\circ}$  instead of the correct  $55^{\circ}$ . As a result, the rotation estimates were considerably smaller than the true values. This also resulted in errors in the translation and depth estimates.

It is important to note that while the magnitude of the angle estimates in [43] were too small, the estimated axis of rotation was correct. In fact, the results show a linear relationship between the true angle and the estimated angle. That the rotation values were always smaller than the true values is consistent with the simulation results shown in Fig. 6. This leads us to the following simple modification to the iterative scheme (Section 4.2):

- 1. For the estimation of motion use a focal length parameter  $f \ll f_{true}$ .
- 2. During the image warping stage use the correct focal length to rotate Images 2 and Image 3 towards Image 1.

This modification required no changes to the structure of the algorithm since the coarse-to-fine and iterative mechanisms were already in place. The results shown in Section 5.4 show that the modified algorithm gives accurate rotation estimates even for medium fields of view ( $55^{\circ}$ ).

4.5.4 How Does This Affect the Translation Estimates? How does using the wrong focal length affect the translation estimates? Looking at (33) (and, also, (29), (32), and (49)), we see that the translations t' and t'' always appear in a dot product with s, where:

$$s = \begin{pmatrix} I_x \\ I_y \\ -xI_x - yI_y \end{pmatrix}.$$

The first and second terms of s are not affected by the focal length estimate. Therefore, there is no effect on the translation estimate in the X and Y direction  $t_x$  and  $t_y$ . In the third term, x and y are supposed to be the normalized image coordinates. If we use the wrong focal length, then the estimated Z translation ( $t_z$ ) will be scaled accordingly.

$$\frac{t_Z}{t_{Ztrue}} = \frac{f}{f_{true}}$$

Using this simple relationship, we can recover the true translation direction from the estimated one even if we use a focal length value which is much smaller than the true value.

# 4.6 Final Touches—Recompute the Depth Keeping the Motion Constant

Initially, during the first iterations and at the coarsest levels of the pyramid, the motion estimates are not accurate and cause errors in the depth estimates. Sometimes, particularly near the borders of the image, these errors are large enough



Fig. 7. The shape of the cost function around the minimum for various fields of view. These simulation results were computed after the program converged to the final motion and depth estimates. For visualization, only two of the 12 parameters were changed at one time. Note that for a narrower field of view (f), the cost function has the shape of a narrow diagonal valley. This indicates an ambiguity between translation along the X axis and rotations around the Y axis.

to affect the new derivatives,  $I'_t$  and  $I''_t$  in such a way that at the next iteration, even when the motion estimates are more accurate, the depth values are still badly in error. This tends to be a local problem and does not affect the global motion estimates but it does create local "holes" in the depth map.

To fix this problem, after computing the motion at the finest level, we go back down the pyramid and recompute the depth while keeping the motion estimate constant. This, in effect, becomes a three camera stereo computation. Fig. 8 shows an example of the depth estimates before and after the second pass.

# 5 EXPERIMENTS AND RESULTS

#### 5.1 Overview

In this section, we perform experiments which test various aspects of the direct estimation of motion and structure. Experiments 5.2 and 5.3 deal with the pure translation case. In Experiment 5.2, we test the accuracy of heading estimation for noncollinear motion. Experiment 5.3 tests a





standard optical flow program on some of the same input images and compares the results with the "direct methods."

We then move on to the case where the motion involves both translation and rotation. Experiment 5.4 revisits the image sequence used in Stein and Shashua [43]. We show that the improvements described in Section 4.5.3 greatly improve the rotation estimates. In Experiment 5.5, we repeat the heading estimation experiments but this time, the motion includes some rotation. We also show that not taking into account even small rotations (~ 0.5°) leads to large errors in heading estimates. In Experiment 5.6, we test the accuracy of Euclidean reconstruction using the simple image of a cube. We measure whether the right-angles between the cube faces are correctly recovered.

In all the preceding experiments, the camera was mounted on a motion stage so that ground truth motion could be known accurately. In Experiment 5.7, we repeat the Euclidean reconstruction experiments, but this time with a hand-held camera. We show that if we neglect the rotation, the shape recovery is noisy. We also show that if we neglect radial lens distortion, then the reconstruction is qualitatively good but the Euclidean measurements, such as angles, are less accurate. This agrees with the observations regarding lens distortion and the trilinear tensor in [41]. Experiments 5.8 and 5.9 show results with simple outdoor scenes.

#### 5.2 Accuracy of Heading Estimates for Pure Translation

In this experiment, we measure the accuracy of the estimation of the camera motion direction (heading direction) in the case of pure translation. We also gauge the effect of the nonlinear lens distortion.



Fig. 9. Schematic diagram of the pure translation experiment. Camera heading is set by the angle  $\alpha$ . Motion 2 and Motion 3 are collinear.

#### 5.2.1 Experimental Procedure

The camera (Pulnix TM9701) was mounted on a motion stage with three degrees of computer controlled motion: horizontal translation, vertical translation, and rotation around the camera's *Y* axis (Fig. 9). The camera lens was (4.9 mm lens /  $82^{\circ}$  FOV).

The camera's *Y* axis was aligned with the vertical axis of the stage. Initially, the camera was positioned so that the optical axis was aligned with the horizontal axis of the translation stage. The accuracy of the alignment was  $\pm 1^{\circ}$ . One can verify the alignment around the *Z* axis by tracking a point as one translates the camera 12 mm in the vertical direction. The *x* coordinate of the point must not move by more than one pixel for every 50 pixels in vertical motion. The alignment between the rotation stage and the vertical translation stage is guaranteed by the stage manufacturer. The alignment around the *X* axis is more difficult to verify but one can compare the parallax due to rotation at point at the bottom and top of the image (see [40]).

The camera was rotated  $\alpha^o$  and an image captured. Then, the camera was translated vertically 12.5 mm and an image captured. Finally, the camera was translated horizontally

10 mm and a third image captured. Thus, we captured an image triplet where the camera motion was pure translation with one of the motions vertical and the second at a heading of  $\alpha^{o}$ . This was repeated for angles  $\alpha = 0^{o}, 10^{o}, \dots, 90^{o}$ .

# 5.2.2 Results

Motion and depth estimation was performed with and without lens distortion correction. Fig. 10 shows the heading estimates. The heading estimates are within a few degrees. Notice that when the true heading was inside the FOV, the errors were less than  $1^{\circ}$ . When the true heading was outside the FOV, the errors increased to  $1 - 2^{\circ}$ . This is a case of a wide angle lens. The errors were considerably larger when lens distortion was not taken into account (solid line in Fig. 10b).

# 5.3 Comparison with Optical Flow Techniques

After computing the motion and depth from an image triplet, we can compute the motion flow using the motion (27). In this experiment, we compare the optical flow estimates obtained in this manner with the flow estimates computed directly from two of the images using an "industry standard" optical flow program. The program is based on code by Bergen and Hingorani of the Sarnoff Corporation [3]. We use two sets of images which have been chosen to be particularly difficult for optical flow programs.

## 5.3.1 Experimental Procedure

Fig. 1 shows three images of a real scene: two paper cylinders, a plaster bust, and a black metal bar in front of a background of vertical stripes. The two motions are pure translation, vertical and horizontal, parallel to the image plane. Fig. 3a shows the recovered depth map for the sequence. The results are qualitatively correct. In Fig. 3b, the texture mapping was removed for clarity. The difficulties in the scene are explained in the Introduction (Section 1) and expanded upon in the results section below.



Fig. 10 (a) Heading estimate from a translating camera. Pure translation was assumed. The  $4.9_{mm}$  lens gives a  $82^{\circ}$  FOV along the X axis. (b) Difference between heading estimate and motion stage reading ("True Heading").



Fig. 11. Optical flow computed for Fig. 1. (a) Correct optical flow computed using depth map and recovered motion. The camera motion was horizontal parallel to the image plane. The flow vectors have a zero Y component. (c) Optical flow computed using code from Bergen and Hingorani. Near the vertical bars on the left, the flow vectors have a strong Y component. This is due to the aperture problem. Even a large window (aperture) will see image gradients in only one direction giving no constraint on the Y component of the flow.

From the depth map and motion, the motion flow was computed using (27). The flow field was also computed using the optical flow program and the results compared.

This experiment was repeated for a second set of images shown in Fig. 14a, 14b, and 14c. The main difficulty here is that most of the gradients are smooth, a low contrast. The strongest edges appear at the occluding contour of the bust and these edges do not represent real features.

### 5.3.2 Results

Fig. 11a shows the estimated flow computed from the recovered depth and motion. The camera motion was horizontal parallel to the image plane. The flow vectors have a zero *Y* component. Fig. 11b shows the magnitude of the flow vectors. Fig. 11c shows the flow estimated from two images using an optical flow program. Near the vertical bars on the left, the flow vectors have a strong *Y* component. This is due to the aperture problem. Even a large window

(aperture) will see image gradients in only one direction giving no constraint on the *Y* component of the flow.

On the cylinder to the lower right of the bust, the magnitude of the flow vectors is too small. This is due to the fact that the edges were horizontal and parallel to the horizontal epipolar lines. This means that the depth estimates for these points are unreliable (Fig. 11d). This problem does not occur in Fig. 11b because two motions, one horizontal and one vertical, were used to compute the original depth map.

An enlarged detail of one of the more difficult regions is shown Fig. 12. A critical error occurs near the intersections (in the image) of the diagonal line and the vertical bars.

The results for the second set of images are shown in Fig. 13. Errors in the flow direction can be seen near the edges of the bust. The most significant errors are errors in magnitude of the flow in the background areas. These can be seen clearly by comparing the mesh plots.



Fig. 12. (a) Detail from the top left corner of Fig. 11a. (b) Detail from the top left corner of Fig. 11c. Along the vertical bars, there is a strong Y component but this is to be expected due to the aperture problem. Lower down, near the intersections (in the image) of the diagonal line and the vertical bars, the flow has a small Y component. This error occurs because the program "tracks" the intersection point as if it were a real "feature" point, but the lines do not intersect in space. Although smaller in magnitude, this is the more significant error because the program will also give a high confidence to this value.

This experiment was intended to highlight the problems with optical flow techniques and how they are overcome by using direct methods. While this was not a simple "straw man" and the optical flow program was one of the best available, a full system based on optical flow would have more components intended to overcome the problems described. Together with flow vectors, the optical flow program would compute a confidence ellipse for each vector. These estimates would then be combined using robust methods to reject outliers during the motion computation. After the camera motion was estimated, the flow would be recomputed using the epipolar constraint and the new flow would be used to compute the depth. Probably three views would be used with different camera motions. Building such a worthy competitor and fine tuning it to these image sets is, of course, beyond the scope of this paper, but this description does give an idea of the complexity of such a system.

# 5.4 Translation and Rotation (Reanalysis of Data from Stein and Shashua, [43])

In this experiment, we test the modified algorithm for recovering both translation and rotation (Section 4.5.3) on the same data set used by Stein and Shashua [43].

### 5.4.1 Experimental Procedure

The images used by [43] were taken with an 8.5 mm lens in the following way: The camera was translated first vertically (10 mm) and then horizontally (5 mm) to the right. At this third position, the camera was rotated to various angles ranging from  $-4.0^{\circ}$  to  $1.0^{\circ}$ . The depth in the scene ranged from 170 mm to 400 mm.

Fig. 14a, 14b, and 14c show three of the input images in the sequence. The flow, due to a rotation of  $-1.8^{\circ}$  in Fig. 14b, is much greater than the flow due to the translation. The direction of the flow, due to rotation by a negative angle, was in the same direction as the image flow induced by the translation.

## 5.4.2 Results

Fig. 14d shows the recovered depth map using images in Fig. 14a, 14b, and 14c. Fig. 14e, and 14f show the 3D rendering of the depth map.

Fig. 15 shows the recovered rotation estimates for true rotations ranging from  $-4.0^{\circ}$  to  $0.8^{\circ}$ . Outside this range, the old version of the algorithm did not converge. Using the new algorithm, the rotation estimates are within 5 percent of the correct value. There is no scaling error as was observed using the old method.

The new algorithm also has a wider range of convergence. The original images were not available for larger rotations, but a similar setup yielded good results for rotation for the range of  $-5.0^{\circ}$  to  $5.0^{\circ}$ . Since rotations can create large amounts of image flow, the key to the successful convergence is having strong signals with low spatial frequency such as the white head on a darker background. This leads to convergence at the coarse level which is then used as the starting point for motion estimates at finer levels. Strong repetitive high frequency patterns, such as checkerboard patterns, are more difficult.

#### 5.5 Motion Heading Estimation with Rotation

The aim of this experiment is to test the motion estimation algorithm over a wide variety of motion directions when the motion included some rotation. For example, this experiment tests whether the algorithm works when the direction of motion is inside and outside the FOV. We also find the error created by neglecting even small amounts of rotation. The two translations are noncollinear.

#### 5.5.1 Experimental Procedure

The camera setup is similar to that of Experiment 1 (Section 5.2). The camera is rotated to a particular heading direction  $(0^o, 10^o, \ldots, 90^o)$ . The camera was translated vertically and then horizontally. After the horizontal translation, the camera was rotated  $\pm 0.5^o$  and  $\pm 1.0^o$ . Therefore, for each image heading, we have an image triplet which is translation only and then triplets where the second motion also includes some small rotations. The rotation axis of the stage passed within 5 mm of the camera center of projection thus rotations of  $1.0^o$  produce translations of less than 1 mm. Fig. 16 shows input images where the horizontal translation heading was  $60^o$  to the camera optical axis.



Fig. 13. Optical flow computed for Fig. 14. (a) Correct optical flow computed using depth map and recovered motion. The camera motion was horizontal parallel to the image plane. The flow vectors have a zero *Y* component. (c) Optical flow computed using code from Bergen and Hingorani. Errors in flow direction can be seen near the edges of the bust. The most significant errors are errors in magnitude of the flow in the background areas. These can been seen clearly by comparing the mesh plots.

#### 5.5.2 Results

Depth and motion estimates were computed either assuming pure translation or allowing possible rotation. The results for various motion and motion-assumption combinations are shown in Fig. 17, Fig. 18, and Fig. 19. We notice that for even a small rotation  $(0.5^{\circ})$ , the heading estimate is off by over  $10^{\circ}$  if we do not take the rotation into account. When we estimate both rotation and translation, the motion estimates are good (RMS error  $2.1^{\circ}$ ) although not as good as those obtained from a purely translating camera under the pure translation assumption (RMS error  $1.44^{\circ}$ ).

Fig. 16e shows the depth map estimated from the images in Fig. 16a, 16c, and 16d allowing for both translation and rotation. The depth map is qualitatively correct even though the second camera motion included  $-1.0^{\circ}$  rotation. Rotations of this order of magnitude are significant. Fig. 16f shows the depth map estimated assuming pure translation when in fact the second camera motion included  $0.5^{\circ}$  rotation.

# 5.6 Euclidean Structure Estimation in the Presence of Rotation

In this section, we evaluate the Euclidean structure estimation when the motion includes rotation. We use a simple scene, a cube, in front of a planar background. We perform the 3D reconstruction of the scene. We then measure whether the right-angles between the cube faces are correctly recovered.

#### 5.6.1 Experimental Procedure

Fig. 20 shows the simple scene of a textured cube in front of a flat, low-contrast, background. The camera (8.5 mm lens /  $52^{o}$  FOV) was translated vertically and then horizontally parallel to the image plane. The camera was then rotated to  $+2.0^{o}$  and  $-1.0^{o}$  from its original heading.

Since we were interested in the 3D structure and not just motion estimates, after estimating both structure and motion, we went back down the pyramid and recomputed the structure while keeping the motion estimates constant



Fig. 14. The three input images (a), (b), and (c) used for 3D reconstruction. Motion  $(a \rightarrow b)$  was a horizontal translation with a  $-1.8^{\circ}$  rotation around the *Y* axis. Motion  $(a \rightarrow c)$  was a vertical translation. This is a challenging set of images because the texture is smooth with low contrast. The strongest edge "features" appear along the occluding contour of the head, but this edge does not, in fact, correspond to a real feature in the world. (d) The estimated depth map. The estimated rotation was  $-1.62^{\circ}$ . There are, of course, errors around the boundary of the image where there is no overlap between the images. (e) and (f) Three-dimensional rendering of depth map in (d).

(Section 5.6). The depth estimates before and after the second pass are shown in Fig. 20d and Fig. 20f, respectively. In Fig. 20d, there is a "hole" on the right of the depth map. This is due to the appearance of a white patch near the

border of the image. As with occlusions, this sudden appearance violates the constant brightness assumption. The situation is fixed by the second pass. In general, points near the edge of the depth map are unreliable.

## 5.6.2 Results

The rotation estimates were  $1.68^{\circ}$  and  $-0.67^{\circ}$  for true rotations of  $+2.0^{\circ}$  and  $-1.0^{\circ}$ , respectively. The 3D rendering of the cube is shown in Fig. 21 and Fig. 22. The wire-frame rendering of the overhead views (Fig. 21b and Fig. 22b) show the recovered angle between the cube faces. The estimated angles were  $95^{\circ}$  and  $86^{\circ}$ , respectively. Fig. 21d shows a side view of the cube. The estimated angle between the top and side face of the cube is  $90^{\circ}$ .

# 5.7 Euclidean Structure Estimation Using a Hand-Held Camera

In the previous experiments, the camera was mounted on a motion stage so that ground truth motion values were available. In this experiment, the camera was hand-held to show that the system can deal with the more natural situation. Simple objects were used to create the scene so that it would be easier to evaluate the shape reconstruction.

#### 5.7.1 Experimental Procedure

A progressive scan camera (Pulnix TMC9701) with a wide angle lens (4.9 mm lens / 82° FOV) was hooked up directly to the frame grabber of an SGI Indy workstation. The cameras help in the hand with no tripod. An image sequence of 30 frames was captured while the camera was moved up and down and side to side, the combination of which produced a circular motion. An effort was made to keep the camera rotation down to a minimum.

From the image sequence, three images were selected which seemed to give two distinct motion directions. Fig. 23a, 23b, and 23c show the three input images of the cube.

#### 5.7.2 Results

The motion and depth was estimated for three different cases. First, the motion and depth was estimated allowing for both rotation and translation and taking into account the radial distortion. In the second case, motion was estimated allowing for both rotation and translation but neglecting



Fig. 15. Estimated rotation as a function of true rotation. (a) Old results: Although there is the correct linear relationship, there also appears to be a significant scale error (note Y axis). (b) New results—the rotation estimates are correct.



Fig. 16. Four of the input images used in the experiment to test motion heading estimation in the presence of rotation. Motion (a) to (d) is vertical motion. Motion (a) to (b) is horizontal motion forward and to the left ( $60^{\circ}$  from the camera optical axis). Motion (a) to (c) has the same translation as (a) to (b) but also with a rotation of  $1.0^{\circ}$ . In this case, the flow, due to rotation of  $1.0^{\circ}$ , cancels out the flow, due to translation in the area of the face and more than cancels out the flow in regions of greater depth. (e) Depth map estimated after computing both translation and rotation. The second camera motion included  $-1.0^{\circ}$  rotation. (f) Depth map estimated assuming pure translation when the second camera motion included  $0.5^{\circ}$  rotation. Notice the depth errors in the background. (Lens 8.5 mm lens /  $52^{\circ}$ FOV).

radial distortion. The resulting depth map is qualitatively the same, but the measured angles between the cubes faces are larger and further from the true value of 90° (Fig. 24c and 24d). In the third case, which assumed pure translation, the resulting depth map again appears to be qualitatively the same, but the 3D renderings (Fig. 24e and 24f) clearly show that the results are more noisy and further from true Euclidean reconstruction.

Errors in depth reconstruction, when we neglect the rotation, are to be expected according to [7]. They show that for small total rotation angles (|w|), the error in depth due to an error in rotation ( $|\delta w|$ ) is given by:

$$\frac{\delta Z}{Z} \approx \left(\frac{f}{|u|}\right) (|\delta w|),$$

where *f* is the focal length and |u| is the magnitude of the image motion. In our case, by assuming pure translation, we get a rotation error of  $1.45^{\circ}$  or 0.025rad (Table 1). The average image motions were on the order of 15pixels and the focal length  $f \approx 180$ . So, we expect:



Fig. 17. Heading estimates from a moving camera. (PT-PT) Both true motion and motion model were pure translation.  $(0.5^o\text{-PT})$  True motion included  $0.5^o$  and motion model was pure translation. Note that even a small rotation, such as  $0.5^o$ , is enough to create large heading estimation errors if rotation is not taken into account. (PT-Rot,  $+1.0^o\text{-Rot}$ ,  $-1.0^o\text{-Rot}$ ) Motion model included rotation Actual rotations were  $0^o$ ,  $1.0^o$ ,  $-1.0^o$ . (8.5<sub>mm</sub> lens / 52<sup>o</sup> FOV).



Fig. 18. Heading estimate error: The difference between heading estimate and motion stage reading. Even in the presence of camera rotation ( $+1.0^{\circ}$ -Rot,  $-1.0^{\circ}$ -Rot) the heading estimation errors are small (RMS error  $2.1^{\circ}$ ) although not as good when obtained from purely translating camera under the pure translation assumption (PT-PT) where the RMS error was  $1.44^{\circ}$ . (8.5mm lens / 52° FOV).



Fig. 19. Rotation estimates obtained from a camera translating in different translation directions with some added rotation. Actual rotations were  $0^{o}, 1.0^{o}, -1.0^{o}$ . (camera with  $8.5_{mm}$  lens /  $52^{o}$  FOV).



Fig. 20. A simple scene used to test 3D Euclidean reconstruction: A cube constructed from cork blocks in front of a flat background. Motion (a) to (b) was vertical. Motion (a) to (c) was a sideways translation with  $-1.0^{\circ}$  rotation. Motion (a) to (e) was the same sideways translation but with  $+2.0^{\circ}$  rotation. (d) Initial depth map from images (a), (b), and (c). Note how the disappearance of a white blob in the background near the right edge of the image creates an error in one region of the depth map. (f) After second pass of depth estimation, the "hole" is fixed (see text).

$$\frac{\delta Z}{Z} \approx 0.3$$

Dutta and Snyder [6] refer to depth from image motion. The errors we observe are not so large, but are definitely noticeable. We also note that even if we neglect the rotation, none of the depth estimates were negative so a *depth is positive* constraint, as suggested by [12], would not help in this example.

#### 5.8 Outdoor Scenes—Church Wall

In this experiment, we tested the algorithm on outdoor scenes. These experiments are important because they show that the photometric constraints can be used in uncontrolled lighting conditions. Fig. 25a, 25b, and 25c show three images of a small part of a church wall. The camera was a Sony Hi8 camcorder with the lens open as wide as possible (FOV 40°). The camera was mounted on a tripod. Moving the camera horizontally  $(a \rightarrow b)$  was quite smooth, but moving the camera up and down included about  $0.5^{\circ}$  rotation. Fig. 25d shows the recovered depth map which can be seen to be qualitatively correct. The Euclidean 3D renderings of the depth map are shown in Fig. 26.













~)

Fig. 21. Three-dimensional rendering of the recovered scene. The camera motion included  $-1.0^{\circ}$  rotation. (a) Overhead view. The estimated camera position is marked by the sphere in bottom right-hand corner. (b) An enlarged view of the corner of the cube. The angle between the faces of the cube is estimated  $95^{\circ}$  instead of  $90^{\circ}$ . (c) Enlarged wire-frame side view of the cube. The estimated angle between the top face of the cube and the right face is  $90^{\circ}$ . (d) Side view of the segmented cube from up and to the side. (e) Another view of the scene showing the foreground/background segmentation. The depth estimates of the back panel are noisy because they are further away from the camera and because the texture does not have strong gradients.

#### 5.9 Outdoor Scene 2—Side Entrance

Fig. 27 shows three views of the side entrance to a building. The camera was a progressive scan camera (Pulnix TMC9701) with a wide angle lens (4.9 mm / FOV 82°). The camera output was recorded on a Sony Hi8 camcorder. The camera was mounted on a tripod and an effort was made to reduce the amount of rotation, but rotation was 0.32°. Fig. 27d shows the recovered depth map. The 3D Euclidean rendering of the depth map is shown in Fig. 28.

#### 6 SUMMARY AND FUTURE WORK

We have presented a general relationship between the spatiotemporal derivatives of three frames and the ego-motion



Fig. 22. Three-dimensional rendering of the recovered cube. The camera motion included  $+2.0^{\circ}$  rotation. The cube has been segmented from the background using depth based foreground/background segmentation. (a) View from up and to the side. (b) Wire-frame rendering of the overhead view. From the overhead view, the estimated angle between the cube faces is  $86^{\circ}$  instead of  $90^{\circ}$ .



Fig. 23. A simple scene used to test 3D Euclidean reconstruction. Camera was hand-held so exact motion between images (a), (b), and (c) is not known. (d) and (e) Three-dimensional renderings of the recovered shape. The motion estimation allowed for rotation and translation and the input images were preprocessed to correct for lens distortion. The estimated angle between cube faces is  $96^{\circ}$  instead of  $90^{\circ}$ .

parameters of the two motions. This relationship was derived first for the general case where the ego-motion model comprises the 27 coefficients of the trilinear tensor and then for the small-motion model of Longuet-Higgins and Prazdny.

These relationships are model-based brightness constraints which provide a linear constraint per pixel in the image—thereby providing a method for direct structure and motion estimation that cuts through the aperture problem and without prior detection of feature points. The linear constraints of the 27-parameter projective model are degenerate and cannot lead directly to a unique solution for the tensor coefficients although a unique solution can be found using quadratic admissibility constraints. This added complexity led us to implement, in practice, the simpler small-motion model of Longuet-Higgins and Prazdny.

The implementation details of these model-based brightness constraints are important and include four critical elements: 1) embedding of the computations within a coarse-to-fine (Gaussian pyramid) framework, 2) Newton





(a)







Fig. 24. Top view of the cube scene in Fig. 23. (a) and (b) Motion and depth was estimated allowing for rotation. The images were preprocessed for radial distortion correction. The estimated angle between cube faces was  $97^{\circ}$ . (c) and (d) No radial distortion correction. The estimated angle between cube faces was  $101^{\circ}$ . (e) and (f) Motion was estimated assuming pure translation. There are many errors in the depth map and the angles estimate is  $105^{\circ}$ .

iterations over the brightness constraint equation, 3) postprocessing of smoothness and surface interpolation for obtaining visually pleasing results, and 4) obtaining both stability and accuracy by assuming a very short focal length for the motion estimation, but using the correct focal length for the image warping.

The algorithm was tested on a set of challenging real image situations which contain very few "good features" (local regions with significant variability of gradient

 TABLE 1

 Rotation Estimates (in Degrees) from the Images in Fig. 23

Ī		$w_x$	$w_y$	$w_z$
	Motion 1	-0.7809	0.2262	-0.5520
	Motion 2	-0.0538	0.2684	1.4246

The camera was hand-held so ground truth values of the motion are not known.



Fig. 25. Detail of church wall: Three input views and resulting depth map. (Sony Hi8 camcorder, FOV  $40^{\circ}$ ).

direction) which are necessary for optical flow and discrete point matching algorithms. Yet, we obtained a faithful dense depth map of the scene due to the fact that the algorithm is not hindered by the presence of aperture effects.

## 6.1 Future Work—Collinear Motion

The method, as described so far, fails when the two motions are in the same (or opposite) directions. This is a drawback in many applications (e.g., 3D reconstruction from a monocular image sequence). Here, we present some early research on ways to overcome this problem. The initial results look promising ([42]). We will investigate the pure translation case because if it fails for pure translation, it will also fail if there is some rotation.

For pure translation:

$$I''_{t}s^{\top}t' - I'_{t}s^{\top}t'' = 0$$
(60)

and

$$ks^{\top}t'' + I''_{t} = 0 \tag{61}$$

$$ks^{\top}t' + I'_t = 0. (62)$$

If the second translation vector t'' is proportional to the translation vector t', then (61) is simply a scaled version of (62) adding no new information and the solution is ill-conditioned. Another way of interpreting (62) and (61) is that a motion t' in one direction will create a change in the image  $I'_t$  which is exactly the opposite of the change  $I''_t$  induced by a motion t'' of equal magnitude and opposite direction to t'.

The LH&P model which assumes  $\binom{t_z}{Z} << 1$  and was used to derive (62) and (61) is of course not exact unless the motion in the *Z* direction is zero, but in any case, the contribution to the optical flow due to translation in the *Z* direction, is small unless the field of view is very wide. Therefore, for small motion, the LH&P model will be quite accurate and the equations will be ill-conditioned. This situation is closely related to the case of reconstruction from line correspondences. There is a critical combination of a line configuration and camera motions which can be stated as follows:





Fig. 26. (a) Three-dimensional Euclidean rendering of the depth map in Fig. 25d. (b) Enlarged view from a slightly different angle.

**Conjecture 1.** Critical configuration of lines and motion. Let S be a set of lines in 3D which have a common intersecting line L (i.e.,  $S \land L = 0$  for all  $S \in S$ ). Let s, s', s'' be the projections of the line  $S \in S$  in the three views: Image 1, Image 2, and Image 3. Let o, o', and o'' be the corresponding camera centers. If o, o', and o'' all lie on a line  $S_i \in S$ , the corresponding lines s, s', s'' do not provide a unique solution to the camera motion.

We currently have no proof, but simulation experiments confirm the theorem. This line configuration is the Linear Line Complex (LLC) which is examined in [44]. In [44], it is shown that for general motion, the tensor has a linear degeneracy in the case of an LLC, but a unique solution is obtained by taking into account nonlinear constraints. Here,



Fig. 27. (a), (b), and (c) Three input views of a side entrance to a building and resulting depth map (d). Camera was mounted on a tripod but motion included up to  $0.5^{\circ}$  rotation.

for the specific case where the camera motion is along a line belonging to the LLC, there is a whole family of solutions. The epipole could lie anywhere on the projection of the line L onto the image. Thus, there is one degree of uncertainty in addition to the scale factor ambiguity.

If we take the common to be the line at infinity on the *XY* plane, we get the particular case where all the lines lie on planes which are parallel to the image plane, and the translations are collinear and also parallel to the image plane. This is a critical condition according to Conjecture 1 and the direction of translation in the *XY* plane cannot be recovered. In such a configuration, lines in Image 1 will be parallel to corresponding lines in Image 2 and Image 3 (assuming no rotation).

How does this relate to our case? By using the first order approximation of the optical flow constraint equation to get the point-line-line correspondences, we have created, as an artifact, a situation where all corresponding lines are parallel. In reality, this is not the case unless all the lines came from planes parallel to the image plane.

# 6.1.1 A Possible Solution

The Optical Flow Constraint Equation [17] was based on the first order Taylor expansion of  $I(x, y, t) = I(x', y', t + \delta t)$ , where *t* is used to denote time (not translation). Keeping the second order terms of the Taylor expansion results in:

$$I(x + \delta x, y + \delta y, t + \delta t) =$$

$$I(x, y, t) + \delta x I_x + \delta y I_y + \delta t I_t +$$

$$\frac{1}{2!} \left( \delta x^2 E_{xx} + \delta y^2 E_{yy} + \delta t^2 E_{tt} + 2\delta x \delta y E_{xy} + 2\delta x \delta t E_{xt} + 2\delta y \delta t E_{yt} \right) + e.$$
(63)

Replacing  $\delta x = u' \delta t$  and  $\delta y = v' \delta t$  yields:



Fig. 28. Three-dimensional renderings of the depth map in Fig. 27d. The spheres in (a) and (b) indicate the estimated camera location. Note in overhead view (d), true Euclidean structure is recovered with correct  $90^{\circ}$  angles.

$$I(x + u\delta t, y + v\delta t, t + \delta t) =$$

$$I(x, y, t) + u'\delta tI_x + v'\delta tI_y + \delta tI_t +$$

$$\frac{1}{2!} \left( (u\delta t)^2 E_{xx} + (v'\delta t)^2 E_{yy} + \delta t^2 E_{tt} +$$

$$2u'v'\delta t^2 E_{xy} + 2u'\delta t^2 E_{xt} + 2v'\delta t^2 E_{yt} \right)$$

$$+ e.$$
(64)

We notice that two of the second order terms are also linear with u and v. If we keep those two terms, (at this point, apart from convenience, we have no formal justification why we should keep those two terms and not all the second order terms) and apply the constant brightness assumption:

$$I(x, y, t) = I(x + \delta x, y + \delta y, t + \delta t),$$

we get for the first motion:

$$u'(I_x + I_{xt}\delta t) + v'(I_y + I_{yt}\delta t) + I'_t = 0.$$
 (65)

Since  $I'_x = I_x + I_{xt}\delta t$  and  $I'_y = I_y + I_{yt}\delta t$  this becomes:

$$u'I'_x + v'I'_y + I'_t = 0. (66)$$

The original work by Horn and Schunk who used:  $u'\frac{(I'_x+I_x)}{2} + v'\frac{(I'_y+I_y)}{2} + I'_t = 0.$ ) Note that  $I'_x$  and  $I'_y$  are computed at coordinates (x, y) not (x', y'). Similarly, for the second motion:

$$u''I''_{x} + v''I''_{y} + I''_{t} = 0. (67)$$

Likewise,  $I''_x, I''_y$  are the brightness gradients in Image 3 at coordinates (x, y). That leads us to the equation:

$$I''_{t}s'^{\top}t' - I'_{t}s''^{\top}t'' = 0, (68)$$

where:

$$s' = \begin{pmatrix} I'_x \\ I'_y \\ -xI'_x - yI'_y \end{pmatrix} \qquad s'' = \begin{pmatrix} I''_x \\ I''_y \\ -xI''_x - yI''_y \end{pmatrix}.$$
 (69)

Now since  $s' \neq s''$ , the equations are better conditioned. This depends on how much change there was in the images" gradients. As we have noted, this depends on whether the surface is parallel to the image plane or not.

A modification is required to the coarse-to-fine framework and the iterative refinement which are described in Section 4.4 and Section 4.2, respectively. For motion estimation, Image 1 is warped (using forward warping) towards Image 2 and Image 3. For depth estimation, Image 2 and Image 3 are warped towards Image 1.

We have shown a solution to the case of collinear motion with pure translation and these ideas have been successfully implemented these ideas. It is not clear how to transfer this to the general motion case and this remains an open area of study.

#### ACKNOWLEDGMENTS

The authors would like to acknowledge the US-IS BSF contract 94-00120 and the European ACTS project AC074. General support for Gideon P. Stein comes from DARPA contracts N00014-94-01-0994 and 95009-5381.

## REFERENCES

- A. Azarbayejani and A.P. Pentland, "Recursive Estimation of Motion, Structure, and Focal Length." *IEEE Trans. Pattern Analysis* and Machine Intelligence, vol. 17, no. 6, pp. 562-575, June 1995.
- [2] J.R. Bergen, P. Anandan, K.J. Hanna, and R. Hingorani, "Hierarchical Model-Based Motion Estimation," *Proc. European Conf. Computer Vision*, June 1992.
- [3] J. Bergen and R. Hingorani, Hierarchical Motion-Based Frame Rate Conversion," technical report, David Sarnoff Research Center, 1990.
- [4] P.J. Burt and E.H. Adelson, "The Laplacian Pyramid as a Compact Image Code," *IEEE Trans. Comm.*, vol. 31, pp. 532-540, 1983.
  [5] R. Deriche and G. Giraudon, "Accurate Corner Detection: An
- [5] R. Deriche and G. Giraudon, "Accurate Corner Detection: An Analytical Study," *Proc. Int'l Conf. Computer Vision*," pp. 66-70, Dec. 1990.
- [6] R. Dutta and M. Snyder, "Robustness of Correspondence-Based Structure from Motion," Proc. Int'l Conf. Computer Vision, pp. 106-110, Dec. 1990.
- [7] O. Faugeras and T. Papadopoulo, "A Nonlinear Method for Estimating the Projective Geometry of Three Views." Proc. Int'l Conf. Computer Vision '98, Jan. 1998.
- [8] O.D. Faugeras, Three Dimensional Computer Vision: A Geometric Viewpoint. Cambridge, Mass.: MIT Press, 1993.
- [9] O.D. Faugeras and B. Mourrain, "On the Geometry and Algebra of the Point and Line Correspondences between N Images," Proc. Int'l Conf. Computer Vision, pp. 951-956, June 1995.
- [10] C. Fermuller and Y. Aloimonos, "Global Rigidity Constraints in Image Displacement Fields," Proc. Int'l Conf. Computer Vision, June 1995.
- [11] C. Fermuller and Y. Aloimonos, "The Confounding of Translation and Rotation in Reconstruction from Multiple Views," Proc. Conf. Computer Vision and Pattern Recognition, June 1997.
- [12] C. Fermuller and Y. Aloimonos, "On the Geometry of Visual Correspondence," *Int'l J. Computer Vision*, vol. 21, no. 3, pp. 223-47, Feb. 1997.
- [13] R. Hartley, "A Linear Method for Reconstruction from Lines and Points," Proc. Int'l Conf. Computer Vision, June 1995.
- [14] D.J. Heeger and A. Jepson, "Simple Method for Computing 3D Motion and Depth," Proc. Int'l Conf. Computer Vision, pp. 96-100, Dec. 1990.
- [15] J. Heel, "Direct Estimation of Structure and Motion from Multiple Frames," AI Memo 1,190, Artificial Intelligence Laboratory, Massachusetts Inst. Technology, Mar. 1990.

- [16] A. Heyden, "Reconstruction from Image Sequences by Means of Relative Depths," Proc. Int'l Conf. Computer Vision, pp. 1,058-1,063, June 1995.
- [17] B.K.P. Horn and B.G. Schunk, "Determining Optical Flow," *Artificial Intelligence*, vol. 17, pp. 185-203, 1981.
- [18] B.K.P. Horn, E.J. Weldon Jr., "Direct Methods for Recovering Motion," Int'l J. Computer Vision, vol. 2, pp. 51-76, 1988.
- [19] M. Irani, B. Rousso, and S. Peleg, "Recovery of Ego-Motion Using Image Stabilization," Proc. Conf. Computer Vision and Pattern Recognition, pp. 454-460, June 1994.
- [20] R. Kumar and P. Anandan, "Direct Recovery of Shape from Multiple Views: A Parallax-Based Approach," Proc. Int'l Conf. Pattern Recognition, Oct. 1994.
- [21] H. Longuet-Higgins and K. Prazdny, "The Interpretation of a Moving Retinal Image," Proc. Royal Soc. London B, vol. 208, pp. 385-397, 1980.
- [22] B. Lucas and T. Kanade, "An Iterative Image Registration Technique with an Application to Stereo Vision," Proc. Int'l Joint Conf. Artifical Intelligence, pp. 674-679, 1981.
- [23] I.S. McQuirk, "An Analog VLSI Chip for Estimating the Focus of Expansion," Artifical Intelligence Technical Report 1,577, Massechusetts Inst. Technology, 1996.
- [24] D. Michaels, "Exploiting Continuity-in-Time in Motion Vision," PhD thesis, Dept. Electrical Eng. and Computer Science, Massachusetts Inst. Technology, May 1992.
- [25] J.J. More, B.S. Garbow, and K.E. Hillstrom, "User Guide for Minpack-1," Technical Report ANL-80-74, Argonne Nat'l Laboratory, Argonne, Ill., Aug. 1980.
- [26] H. . Nagel, "Direct Estimation of Optical Flow and of Its Derivates," Artificial and Biological Vision Systems, G.A. Orban and H.H. Nagel, eds., pp. 193-224, Springer-Verlag, 1992.
- [27] S. Negahdaripour and B.K.P. Horn, "Direct Passive Navigation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 9, no. 1, pp. 168-176, 1987.
- [28] W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge Univ. Press, second ed., 1992.
- [29] J. Reiger and D.T. Lawton, "Processing Differential Image Motion," J. Optical Soc. Am., vol. 2, pp. 354-359, 1985.
- [30] H. Sawhney, "3D Geometry from Planar Parallax," Proc. Conf. Computer Vision and Pattern Recognition, pp. 929-934, June 1994.
- [31] A. Shashua, "Algebraic Functions for Recognition"IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 17, no. 8, pp. 779-789, Aug. 1995.
- [32] A. Shashua and S. Avidan, "The Rank4 Constraint in Multiple View Geometry," Proc. European Conf. Computer Vision, Apr. 1996.
- [33] A. Shashua and N. Navab, "Relative Affine Structure: Theory and Application to 3D Reconstruction from Perspective Views," Proc. Conf. Computer Vision and Pattern Recognition, pp. 483-489, 1994.
- [34] A. Shashua and N. Navab, "Relative Affine Structure: Canonical Model for 3D from 2D Geometry and Applications," *IEEE Trans. Pattern Analysis and Machine Inteliigence*, vol. 18, no. 9, pp. 873-883, Sept. 1996.
- [35] A. Shashua and P. Anandan, "Trilinear Constraints Revisited: Generalized Trilinear Constraints and the Tensor Brightness Constraint," *Proc. ARPA Image Understanding Workshop*, Feb. 1996.
- [36] A. Shashua and K.J. Hanna, "The Tensor Brightness Constraints: Direct Estimation of Motion Revisited," technical report, Technion, Haifa, Israel, Nov. 1995.
- 37] A. Shashua and M. Werman, "Trilinearity of Three Perspective Views and Its Associated Tensor," *Proc. Int'l Conf. Computer Vision*, pp. 920-925, June 1995.
- [38] S. Soatto and P. Perona, "Motion from Fixation," Proc. Conf. Computer Vision and Pattern Recognition, pp. 817-824, June 1996.
- [39] M. Spetsakis and J. Aloimonos, "A Unified Theory of Structure from Motion," Proc. ARPA Image Understanding Workshop, 1990.
- [40] G.P. Stein, "Internal Camera Calibration Using Rotation and Geometric Shapes," master's thesis, Massachusetts Inst. Technology, Cambridge, Mass., 1993.
- [41] G.P. Stein, "Lens Distortion Calibration Using Point Correspondences," Proc. Conf. Computer Vision and Pattern Recognition, June 1997.
- [42] G.P. Stein, "Geometric and Photometric Constraints: Motion and Structure from Three Views," PhD thesis, MIT Artificial Intelligence Laboratory, Feb. 1998.

1014

- [43] G.P. Stein and A. Shashua, "Model-Based Brightness Constraints: On Direct Estimation of Structure and Motion" Proc. Conf. Computer Vision and Pattern Recognition, June 1997.
- [44] G.P. Stein and A. Shashua, "On Degeneracy of Linear Reconstruction from Three Views: Linear Line Complex and Applications," *Proc. European Conf. Computer Vision*, June 1998.
- [45] R. Szeliski and S.B. Kang, "Direct Methods for Visual Scene Reconstruction," Proc. IEEE Workshop Representation of Visual Scenes, June 1995.
- [46] T.Y. Tian, C. Tomasi, and D.J. Heeger, "Comparison of Approaches to Egomotion Computation," *Proc. Conf. Computer Vision and Pattern Recognition*, pp. 315-320, June 1996.
- [47] C. Tomasi and T. Kanade, Shape and Motion from Image Streams—A Factorization Method," Int'l J. Computer Vision, vol. 9, no. 2, pp. 137-154, 1992.
- [48] B. Triggs, "Matching Constraints and the Joint Image," Proc. Int'l Conf. Computer Vision, pp. 338-343, June 1995.
- [49] J. Weng, N. Ahuja, and T. Huang, "Two View Matching," Proc. Conf. Computer Vision and Pattern Recognition, pp. 64-73, Dec. 1989.
- [50] G. Wolberg, Digital Image Warping. Los Alamitos, Calif.: IEEE Computer Soc. Press, 1990.



Gideon P. Stein received the BS degree in 1990 from the Technion, Haifa, Israel, the MSc degree (May 1993), and the PhD degree (June 1998) from the Artificial Intelligence Laboratory, MIT, Cambridge, Massachusetts. His PhD thesis was titled: "Geometric and Photometric Constraints: Motion and Structure from Three Views." From 1998 to 1999, he was a postdoctoral fellow at the MIT Artificial Intelligence Laboratory working on the Visual Surveillance

and Monitoring project. Since 1999, he has been the vice president of Research and Development at MobilEye Vision Technologies Ltd. in Jerusalem, Israel. Mobileye is a young company focusing on developing computer vision based systems for driving assistance and enhanced road safety.



Amnon Shashua received the BSc. degree in mathematics and computer science from Tel-Aviv, Israel, in 1986, the MSc degree in mathematics and computer science from the Weizman Institute of Science, Rehovot, Israel, in 1989, and the PhD degree in computational neuroscience, working at the Artificial Intelligence Laboratory, from the Massachusetts Institute of Technology, in 1993. He is currently a senior lecturer at the Institute of Computer Science,

The Hebrew University of Jerusalem. His research interests are in computer vision and computational modeling of human vision. His previous work includes early visual processing of saliency and grouping mechanisms, visual recognition, image synthesis for animation and graphics, and theory of computer vision in the areas of threedimensional processing from a collection of two-dimensional views. He is a member of the IEEE.