

Linear and Incremental Acquisition of Invariant Shape Models from Image Sequences

Daphna Weinshall and Carlo Tomasi

Abstract—

We show how to automatically acquire a Euclidean shape representations of objects from noisy image sequences under weak perspective. The proposed method is linear and incremental, requiring no more than pseudo-inverse. A nonlinear but numerically sound preprocessing stage is added to improve the accuracy of the results even further. Experiments show that attention to noise and computational techniques improve the shape results substantially with respect to previous methods proposed for ideal images.

Keywords— Structure from motion, linear reconstruction, factorization method, affine shape, Euclidean shape, weak perspective, Gramian, affine coordinates.

1 Introduction

In model-based recognition, images are matched against stored libraries of three-dimensional object representations, so that a good match implies recognition of the object. The recognition process is greatly simplified if the quality of the match can be determined without camera calibration, namely, without having to compute the pose of each candidate object in the reference system of the camera. For this purpose, three-dimensional object representations have been proposed [15] that are invariant with respect to *similarity* transformations, that is, rotations, translations, and isotropic scaling. These are exactly the transformations that occur in the *weak perspective* projection model, where images are scaled orthographic projections of rotated and translated objects. Because of its linearity, weak perspective strikes a good balance between mathematical tractability and model generality.

In this paper, we propose a method for *acquiring* a Euclidean representation from a sequence of images of the objects themselves. Automatic acquisition from images avoids the tedious and error prone process of typing three-dimensional coordinates of points on the objects, and makes expensive three-dimensional sensors such as laser rangefinders unnecessary. However, model recognition techniques such as geometric hashing have been shown [2] to produce false positive matches with even moderate levels of error in the representations or in the images. Consequently, we pay close attention to accuracy and numerical soundness of the algorithms employed, and derive a computationally robust and efficient counterpart to the schemes that previous

papers discuss under ideal circumstances.

To be sure, several systems have been proposed for computing depth or shape information from image sequences. For instance, [14, 9] identify the minimum number of points necessary to recover motion and structure from two or three frames, [1] recovers depth from many frames when motion is known, [8] considers restricted or partially known motion, [12] solves the complete multiframe problem under orthographic projection, and [5] proposes multiframe solutions under perspective projection.

Conceivably, one could use one of these algorithms to determine the complete three-dimensional shape and pose of the object in a Euclidean reference system, and process the results to achieve similarity invariance. However, a Euclidean representation is *weaker* than a full representation with pose, since it does not include the orientation of the camera relative to the object. Consequently, the invariant representation contains less information, and ought to be easier to compute. This intuition is supported by experiments with complete calibration and reconstruction algorithms, which, given a good initial guess of the shape of the object, spend a large number of iterations modifying the parameters of the calibration and pose matrices, without affecting the shape by much¹.

In this paper we show that this is indeed the case. (We assume weak perspective projection.) Specifically, we compute a similarity-invariant (Euclidean) representation of shape both *linearly* and *incrementally* from a sequence of weak perspective images. This is a very important gain. In fact, a linear multiframe algorithm avoids both the instability of two- or three-frame recovery methods and the danger of local minima that nonlinear multiframe methods must face. Moreover, the incremental nature of our method makes it possible to process images one at a time, moving away from the storage-intensive batch methods of the past.

Our acquisition method is based on the observation that the trajectories that points on the object form in weak perspective image sequences can be written as linear combinations of three of the trajectories themselves, and that the coefficients of the linear combinations represent shape in an affine-invariant basis. This result is closely related to, but different from, the statement that any image in the sequence is a linear combination of three of its images [13].

In this paper, we also show that the optional addition of a nonlinear but numerically sound stage, which selects the

D. Weinshall is with the Institute of Computer Science, Hebrew University of Jerusalem, 91904 Jerusalem, Israel; email: daphna@cs.huji.ac.il. C. Tomasi is with the Department of Computer Science, Stanford University, Cedar Hall, Stanford, CA 94305; email: tomasi@cs.stanford.edu. IEEECS Log Number P95063.

¹B. Boufama, personal communication.

most suitable basis trajectories, improves the accuracy of the representation even further. This leads to an image-to-model matching criterion that better discriminates between new images that depict the model object and those that do not. In order to compare our method to existing model acquisition (or structure from motion) methods, we describe a simple transformation, by which we compute a depth representation from the Euclidean representation computed by our algorithm.

In the following, we first define the weak perspective imaging model (Section 2). We review the Euclidean shape representation and the image-to-model matching measure (Section 3). We then introduce our linear and incremental acquisition algorithm, as well as the nonlinear preprocessing procedure (Section 4). Finally, we evaluate performance with some experiments on real image sequences (Section 5).

2 Multiframe Weak Perspective

Under weak perspective, a point $\mathbf{p}_n = (X_n, Y_n, Z_n)^T$ on an object can be related to the corresponding image point $\mathbf{w}_{mn} = (\xi_{mn}, \eta_{mn})^T$ in frame m by a scaling, a rotation, a translation, and a projection:

$$\mathbf{w}_{mn} = \Pi(s_m R_m \mathbf{p}_n + \mathbf{t}_m) \quad (1)$$

where R_m is an orthonormal 3×3 matrix, \mathbf{t}_m is a three-dimensional translation vector, s_m is a scalar, and Π is the orthographic projection operator that simply selects the first two rows of its argument. The two components of \mathbf{w}_{mn} are thus:

$$\xi_{mn} = s_m \mathbf{i}_m^T \mathbf{p}_n + a_m, \quad \eta_{mn} = s_m \mathbf{j}_m^T \mathbf{p}_n + b_m \quad (2)$$

where the orthonormal vectors $\mathbf{i}_m^T, \mathbf{j}_m^T$ are the first two rows of R_m , and a_m, b_m are the first two components of \mathbf{t}_m . In a sequence of images, feature points can be extracted and tracked (see, e.g., [11]). If N points are tracked in M frames, the equations (2) are repeated MN times, and can be written in matrix form as follows:

$$\begin{bmatrix} \xi_{11} & \dots & \xi_{1N} \\ \vdots & & \vdots \\ \xi_{M1} & \dots & \xi_{MN} \\ \eta_{11} & \dots & \eta_{1N} \\ \vdots & & \vdots \\ \eta_{M1} & \dots & \eta_{MN} \end{bmatrix} = \begin{bmatrix} s_1 \mathbf{i}_1^T \\ \vdots \\ s_M \mathbf{i}_M^T \\ s_1 \mathbf{j}_1^T \\ \vdots \\ s_M \mathbf{j}_M^T \end{bmatrix} \begin{bmatrix} \mathbf{p}_1 & \dots & \mathbf{p}_N \end{bmatrix} + \begin{bmatrix} a_1 \\ \vdots \\ a_M \\ b_1 \\ \vdots \\ b_M \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}$$

that is,

$$\hat{W} = RP + \mathbf{t}\mathbf{1}^T \quad (3)$$

where $\mathbf{1}$ is a vector of N ones. Thus, \hat{W} collects the image measurements, R represents both scaling and rotation in the M frames, P is shape, and \mathbf{t} is translation. In Section 4, we show that R and P need not in fact be computed explicitly in order to compute a Euclidean representation.

3 Review of the Euclidean Representation

Starting with Eq. (3) as a multiframe imaging model, we now describe how to define a shape representation that is invariant with respect to similarity transformations, that is, rigid transformations and isotropic scaling [15]. Specifically, we work towards similarity invariance in three steps:

1. invariance to translation: we use the centroid of the points as a reference origin in the coordinate system where P is described. We translate \hat{W} accordingly, obtaining the matrix of centered image measurements W . Eq. (3) becomes

$$W = RP. \quad (4)$$

2. invariance to affine transformations (Section 3.1);
3. invariance to similarity transformations (Section 3.2).

For performance evaluation only, we will also discuss the

4. computation of depth (Section 3.3).

3.1 Affine Transformation Invariance

The $M \times 3$ matrix R in Eq. (4) is built from 3×3 orthonormal matrices and isotropic scaling factors (see Eq. (2)). Therefore corresponding rows in the upper and lower halves of R (that is, rows m and $m + M$ for $m = 1, \dots, M$) must be mutually orthogonal and have the same norm s_m . If these *orthogonality* constraints are satisfied, we say that R, \mathbf{t} represent full Euclidean motion, and the corresponding P represents Euclidean shape. In particular, the columns of P are the three-dimensional coordinates of the object points with respect to some orthonormal reference basis.

Invariance with respect to affine transformations is achieved by replacing this basis by one that is more intimately related to the shape of the object. Specifically, the basis is made by three of the object points themselves, that is, by the vectors from the reference origin to the three points, assumed not to be coplanar with the origin. This basis is no more orthonormal. The new coordinates were called *affine* in [7]. If now the object undergoes some affine transformation, so do the basis points, and the affine coordinates of the N object points do not change.

The choice of the three basis points can be important. In fact, the requirement that the points be noncoplanar with the origin is not an all-or-nothing proposition. Four points can be *almost* coplanar, and with noisy data this is almost as bad as having *exactly* coplanar points. We discuss this issue in Section 4, where we propose a method that selects a basis as far away as possible from being coplanar with the origin.

Notice that in the new affine basis the three selected basis points have coordinates $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$, so that the new $3 \times N$ matrix A of affine coordinates is related to the Euclidean matrix P of Eq. (4) by the 3×3 linear transformation:

$$P_b A = P \quad (5)$$

where $P_b = [\mathbf{p}_i \ \mathbf{p}_j \ \mathbf{p}_k]$ is the submatrix of P that collects the three selected basis points.

If we substitute Eq. (5) into Eq. (4), we obtain

$$W = W_b A \quad (6)$$

where $W_b = R P_b$. However, because the submatrix $A_b = [\mathbf{a}_i \ \mathbf{a}_j \ \mathbf{a}_k]$ is the identity matrix, we see that W_b is a submatrix of W :

$$W_b = \begin{bmatrix} \mathbf{w}_i & \mathbf{w}_j & \mathbf{w}_k \end{bmatrix}.$$

In more geometric terms, Eq. (6) expresses the following key result:

all the image trajectories (W) of the object points can be written as a linear combination of the image trajectories (W_b) of three of the points. The coefficients (A) of the linear combinations are the three-dimensional coordinates of the corresponding points in space in the affine three-dimensional basis of the points themselves.

Notice the analogy and difference between this result and the statement, made in [13], that under weak perspective any image of an object is a linear combination of three of its views. We are saying that any trajectory is a linear combination of three trajectories, while they are saying that any snapshot is a linear combination of three snapshots. The concise matrix equation in (6) contains these two statements in a symmetric form: Ullman and Basri read the equation by rows, we read it by columns.

3.2 Similarity Invariance

To achieve invariance with respect to similarity transformations, we augment the affine representation introduced above with metric information about the three basis points. Of course, we cannot simply list the coordinates of the three basis points in a fixed reference system, since these coordinates would not be invariant with respect to rotation and scaling. Instead, we introduce the *Gramian* matrix of the three basis points, defined as follows [15]:

$$G = \begin{bmatrix} \mathbf{p}_i^T \mathbf{p}_i & \mathbf{p}_i^T \mathbf{p}_j & \mathbf{p}_i^T \mathbf{p}_k \\ \mathbf{p}_j^T \mathbf{p}_i & \mathbf{p}_j^T \mathbf{p}_j & \mathbf{p}_j^T \mathbf{p}_k \\ \mathbf{p}_k^T \mathbf{p}_i & \mathbf{p}_k^T \mathbf{p}_j & \mathbf{p}_k^T \mathbf{p}_k \end{bmatrix}. \quad (7)$$

In Section 4.2 we normalize G to make it invariant to scaling.

The Gramian is a symmetric matrix, and is defined in terms of the *Euclidean* coordinates of the basis points. However, we show in Section 4 that G can be computed *linearly* from the images, without first computing the depth or pose of the object.

The pair of matrices (A, G) is our target representation. We next show constructively that the pair (A, G) contains complete information about the object's shape, but not directly about its pose in each image.

3.3 Depth Map

Determining the depth of the object requires to express its shape in an orthonormal system of reference, that is, to compute the matrix P of Eq. (4). We now show that the shape Gramian G of Eq. (7) contains all the necessary information. In fact, let W_b be the matrix of the basis trajectories introduced in Eq. (6), and let P_b be the coordinates of the corresponding basis points in space in an orthonormal reference system (see Eq. (5)). Then, the definition (7) of the Gramian can be rewritten as

$$G = P_b^T P_b. \quad (8)$$

Suppose now that T is the Cholesky factor of the Gramian G . We recall that the Cholesky factor of a symmetric positive definite matrix G is the unique upper triangular matrix T with positive diagonal entries such that

$$G = T^T T. \quad (9)$$

Eq. (8) and Eq. (9) are formally similar factorizations of G . We claim that P_b can differ from T only by a rotation or a mirror transformation, so that T is in fact the representation of the three selected basis points on the object in an orthonormal frame of reference. The projection equation (4) does not specify the particular orientation of the orthonormal axes of the underlying reference system, so P_b and T can be taken to be the same matrix up to mirror transformation: $P_b = T$.

In summary, we have the following method for computing the shape matrix P of Eq. (4): determine the Gramian G by the linear method of Section 4, take its Cholesky factorization $T = P_b$, and let T be the transformation of the affine shape matrix A into the new orthonormal basis. Namely:

$$P = T A$$

Notice that the three basis points i, j, k , whose coordinates in A are the identity matrix, are transformed into the columns of T .

We emphasize once more that this last decomposition stage need not be performed for the computation of the Euclidean shape representation. Furthermore, this stage can fail in the presence of noise. In fact, the matrix G can be Cholesky-decomposed only if it is positive definite. Bad data can cause this condition to be violated.

3.4 Relations governing the representations

Once the Euclidean representation (A, G) has been determined from a given sequence of images, it can be used on new, unfamiliar views to determine whether they contain the object represented by (A, G) . In fact, from Eq. (4) and Eq. (8) we obtain

$$W_b G^{-1} W_b^T = R P_b (P_b^T P_b)^{-1} P_b^T R^T = R R^T.$$

If we write out the relevant terms of this equation, we have $\forall m$:

$$\begin{aligned} \mathbf{x}_m^T G^{-1} \mathbf{x}_m &= \mathbf{y}_m^T G^{-1} \mathbf{y}_m = s_m^2 \\ \mathbf{x}_m^T G^{-1} \mathbf{y}_m &= 0 \end{aligned} \quad (10)$$

where the vectors $\mathbf{x}_m^T = (x_{m1}, x_{m2}, x_{m3})$ and $\mathbf{y}_m^T = (y_{m1}, y_{m2}, y_{m3})$ are the rows of the upper and lower half of the centered image measurement matrix W_b . Namely, \mathbf{x}_m and \mathbf{y}_m are the centered image measurement of the basis points in frame m .

Eq. (10) provides strong constraints, capturing all the information that can be obtained from a single image, since all the images that satisfy Eq. (10) are a possible instance of the object represented by G . The two equations in Eq. (10) can be used in two ways: during recognition, the given G can be used to check whether new image measurements \mathbf{x}_m^T , \mathbf{y}_m^T represent the same three basis points as in the familiar views, thus yielding a key for indexing into the object library. During acquisition of the shape representation, on the other hand, G is the unknown, and Eq. (10) can be solved for G .

4 The Algorithm

In this section, we show how to compute the affine shape matrix A (Section 4.1) and the Gramian G of the basis points (Section 4.2) linearly and incrementally from a sequence of images. We then show how to choose three good basis points i, j, k (Section 4.3). This algorithm can use as little as two frames and five points for computing matrix A , and as little as three frames and four points for computing matrix G . More data can be added to the computation incrementally, if and when available.

4.1 The Affine Shape Matrix

The affine shape matrix A is easily computed as the solution of the overconstrained linear system (6), which we repeat for convenience: $W = W_b A$. Recall that W is the matrix of centered image measurements, and W_b is the matrix of centered image measurements of the basis points.

It is well known from the literature of Kalman filtering that linear systems can be solved incrementally one row at a time. The idea is to realize that the expression for the solution

$$A = W_b^+ W$$

where W_b^+ is the pseudoinverse of W_b :

$$W_b^+ = (W_b^T W_b)^{-1} W_b^T$$

is composed of two parts whose size is independent of the number of image frames, namely, the so-called covariance matrix

$$Q = (W_b^T W_b)^{-1}$$

of size 3×3 , and the $3 \times M$ matrix

$$S = W_b^T W.$$

Both Q and S can be updated incrementally every time a new row \mathbf{w}^T is added to W (so the corresponding row \mathbf{w}_b^T is also added to W_b). Specifically, the matrices Q_+ and

S_+ after the update are given by

$$\begin{aligned} Q_+ &= \left(I - \frac{Q \mathbf{w}_b \mathbf{w}_b^T}{1 + \mathbf{w}_b^T Q \mathbf{w}_b} \right) Q \\ S_+ &= S + \mathbf{w}_b \mathbf{w}^T \end{aligned}$$

where I is the 3×3 identity matrix. For added efficiency, this pair of equations can be manipulated into the following update rule for A :

$$A_+ = A + \frac{Q \mathbf{w}_b}{1 + \mathbf{w}_b^T Q \mathbf{w}_b} (\mathbf{w}^T - \mathbf{w}_b^T A).$$

Note that the computation of A requires at least two frames.

4.2 The Gramian

For each frame m , the two equations in (10) define linear constraints on the entries of the inverse Gramian $H = G^{-1}$, so H can be computed as the solution of a linear system. This system, however, is homogeneous, so H can only be computed up to a scale factor.

To write this linear system in the more familiar form $C\mathbf{h} = 0$, we first notice that H is a symmetric 3×3 matrix, so it has six distinct entries h_{ij} , $1 \leq i \leq j \leq 3$. Let us gather those entries in the vector

$$\mathbf{h} = [h_{11} \ h_{12} \ h_{13} \ h_{22} \ h_{23} \ h_{33}]^T.$$

Furthermore, given two 3-vectors \mathbf{a} and \mathbf{b} , define the operator

$$\mathbf{z}^T(\mathbf{a}, \mathbf{b}) = [a_1 b_1 \ a_1 b_2 + a_2 b_1 \ a_1 b_3 + a_3 b_1 \ a_2 b_2 \ a_2 b_3 + a_3 b_2 \ a_3 b_3].$$

Then, the equations in (10) are readily verified to be equivalent to the $2M \times 6$ system

$$C\mathbf{h} = 0 \tag{11}$$

where

$$C = \begin{bmatrix} \mathbf{z}^T(\mathbf{x}_1, \mathbf{x}_1) - \mathbf{z}^T(\mathbf{y}_1, \mathbf{y}_1) \\ \vdots \\ \mathbf{z}^T(\mathbf{x}_M, \mathbf{x}_M) - \mathbf{z}^T(\mathbf{y}_M, \mathbf{y}_M) \\ \mathbf{z}^T(\mathbf{x}_1, \mathbf{y}_1) \\ \vdots \\ \mathbf{z}^T(\mathbf{x}_M, \mathbf{y}_M) \end{bmatrix}.$$

A unit norm solution to this linear system is reliably and efficiently obtained from the singular value decomposition of $C = U_C \Sigma_C V_C^T$ as

$$\mathbf{h} = \mathbf{v}_{C6},$$

the sixth column of V_C . Because this linear system is overconstrained as soon as $M \geq 3$, the computation of H , and therefore of the Gramian $G = H^{-1}$, can be made insensitive to noise if sufficiently many frames are used. Notice that the fact that the vector \mathbf{h} has unit norm automatically normalizes the Gramian.

Alternatively, in order to obtain an incremental algorithm for the computation of the Gramian G , Eq. (11) can be solved with pseudo-inverse. (The incremental implementation of pseudo-inverse was discussed in Section 4.1.) However, the method of choice for solving homogeneous linear systems, which avoids rare singularities, is the method outlined above using SVD.

4.3 Selecting a Good Basis

The computation of the Euclidean representation (A, G) is now complete. However, no criterion has yet been given to select the three basis points i, j, k . The only requirement so far has been that the selected points should not be coplanar with the origin. However, a basis can be very close to coplanar without being strictly coplanar, and in the presence of noise this is almost equally troublesome.

To make this observation more quantitative, we define a basis to be good if for any vector \mathbf{v} the coordinates \mathbf{a} in that basis do not change much when the basis is slightly perturbed. Quantitatively, we can measure the quality of the basis by the norm of the largest perturbation of \mathbf{a} that is obtained as \mathbf{v} ranges over all unit-norm vectors. The size of this largest perturbation turns out to be equal to the *condition number* of W_b , that is, to the ratio between its largest and smallest singular values.

The problem of selecting three columns W_b of W that are as good as possible in this sense is known as the *subset selection problem* in the numerical analysis literature. In the following, we summarize the standard solution to this problem:

1. compute the singular value decomposition of W , $W = U\Sigma V^T$
2. apply QR factorization with column pivoting to the right factor V^T , $V^T = \hat{Q}\hat{R}\Pi^T$

The first three columns of the permutation matrix Π are all zero, except for one entry in each column, which is equal to one. The row subscripts of those three nonzero entries are the desired subscripts i, j, k .

The rationale of this procedure is that singular value decomposition preconditions the shape matrix, and then QR factorization with column pivoting brings a well conditioned submatrix in front of $\hat{Q}\hat{R}$.

Although heuristic in nature, this procedure has proven to work well in all the cases we considered (see analysis of real sequences in [16]). Both the singular value decomposition and the QR factorization of a $M \times N$ matrix can be performed in time $O(MN^2)$, so this heuristical algorithm is much more efficient than the $O(MN^3)$ brute-force approach of computing the condition numbers of all the possible bases.

4.4 Summary of the Algorithm

The following steps summarize the algorithm for the acquisition of the Euclidean representation (A, G) from a sequence \hat{W} of images under weak perspective (see Eq. (3)).

1. Center the measurement matrix with respect to one of its columns or the centroid of all its columns:

$$W = \hat{W} - \mathbf{t}\mathbf{1}^T$$

where \mathbf{t} is either the first column of W or the average of all its columns.

2. (optional) Find a good basis i, j, k for the columns of W as follows:
 - (a) compute the singular value decomposition of W , $W = U\Sigma V^T$
 - (b) apply QR factorization with column pivoting to the right factor V^T , $V^T = \hat{Q}\hat{R}\Pi^T$

The row subscripts of the three nonzero entries in the first three columns of Π are i, j, k . These are the indices of the chosen basis points.

3. Compute the solution A to the overconstrained system $W = W_b A$ by adding one row at a time. Specifically, initialize A to a $3 \times N$ matrix of zeros. Let \mathbf{w}^T be a new row, let \mathbf{w}_b^T collect entries i, j, k of \mathbf{w} , and let $Q = (W_b^T W_b)^{-1}$. The matrix A is updated to

$$A_+ = A + \frac{Q\mathbf{w}_b}{1 + \mathbf{w}_b^T Q \mathbf{w}_b} (\mathbf{w}^T - \mathbf{w}_b^T A) .$$

4. Determine the Gramian G as follows:

- (a) construct the $2M \times 6$ matrix

$$C = \begin{bmatrix} \mathbf{z}^T(\mathbf{x}_1, \mathbf{x}_1) - \mathbf{z}^T(\mathbf{y}_1, \mathbf{y}_1) \\ \vdots \\ \mathbf{z}^T(\mathbf{x}_M, \mathbf{x}_M) - \mathbf{z}^T(\mathbf{y}_M, \mathbf{y}_M) \\ \mathbf{z}^T(\mathbf{x}_1, \mathbf{y}_1) \\ \vdots \\ \mathbf{z}^T(\mathbf{x}_M, \mathbf{y}_M) \end{bmatrix}$$

where

$$\mathbf{z}^T(\mathbf{a}, \mathbf{b}) = [a_1 b_1 \quad a_1 b_2 + a_2 b_1 \quad a_1 b_3 + a_3 b_1 \quad a_2 b_2 \quad a_2 b_3 + a_3 b_2 \quad a_3 b_3] ;$$

- (b) solve the system

$$C\mathbf{h} = 0$$

which yields the distinct entries of the symmetric matrix H . Compute G as the inverse of H .

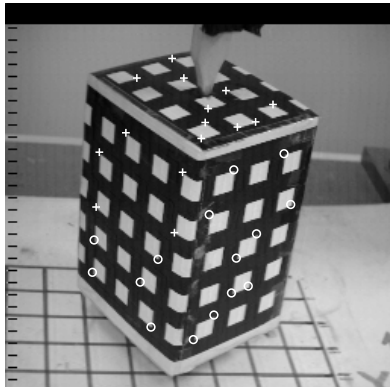
In order to compute a depth map P from the Euclidean representation (A, G) , another optional step is added to the algorithm:

5. (optional) Take the Cholesky factorization of matrix $G = T^T T$, and let T be the transformation of the affine shape matrix A into an orthonormal basis:

$$P = T A$$

5 Experiments

We applied our algorithm, including the depth computation, to two sequences of images, originally taken by Rakesh Kumar and Harpreet Singh Sawhney at UMASS-Amherst (see Fig. 1). The data was provided by J. Inigo Thomas from UMass, who also provided the solution to the correspondence problem (namely, a list of the coordinates of the tracked points in all the frames).



(a)



(b)

Figure 1: (a) One frame from the box sequence, (b) one frame from the room sequence.

For comparison, we received the 3D coordinates of the points in the first frame as ground truth. We used the algorithm described in [3] to compute the optimal similarity transformation between the invariant depth map representation computed by our algorithm (step 5), and the given data in the coordinate system of the first frame. We applied the transformation to our depth reconstruction to obtain z_{est} at each point, and compared this output with the ground truth data z_{real} . We report the relative error at each point, namely, $\frac{z_{est} - z_{real}}{z_{real}}$.

We evaluated the affine shape reconstruction separately. We computed the optimal affine transformation between the invariant affine representation computed by our algorithm (matrix A computed in step 3), and the given depth data in the coordinate system of the first frame. We applied the transformation to the affine shape representation

to obtain z_{est}^{aff} at each point, and compared this output with the ground truth data z_{real} .

5.1 Box sequence:

This sequence includes 8 images of a rectangular chequered box rotating around a fixed axis (one frame is shown in Fig. 1a). 40 corner-like points on the box were tracked. The depth values of the points in the first frame ranged from 550 to 700 mms, therefore weak perspective provided a good approximation to this sequence. (See a more detailed description of the sequence in [10] Fig. 5, or [6] Fig. 2.)

We compared the relative errors of our algorithm to the errors reported in [10]. Three results were reported in [10] and copied to Table 1: column “Rot.” – depth computation with their algorithm, which assumes perspective projection and rotational motion only; column “2-frm” – depth computation using the algorithm described in [4], which uses 2-frames only; and column “2-frm, Ave.” – depth computation using the 2-frames algorithm, where the depth estimates were averages over six pairs of frames. Table 1 summarizes these results, as well as the results using our affine algorithm (column “Aff. Invar.”) and similarity algorithm (column “Rigid Invar.”).

5.2 Room sequence

This sequence, which was used in the 1991 motion workshop, includes 16 images of a robotic laboratory, obtained by rotating a robot arm 120° (one frame is shown in Fig. 1b). 32 corner-like points were tracked. The depth values of the points in the first frame ranged from 13 to 33 feet, therefore weak perspective does *not* provide a good approximation to this sequence. Moreover, a wide-lens camera was used, causing distortions at the periphery which were not compensated for. (See a more detailed description in [10] Fig. 4, or [6] Fig. 3.)

Table 2 summarizes the results of our invariant algorithm for the last 8 points. Due to the noise in the data and the large perspective distortions, not all the frames were consistent with rigid motion. (Namely, when all the frames were used, the computed Gramian was not positive-definite). We therefore used only the last 8 frames from the available 16 frames.

We compared in Table 3 the average relative error of the results of our algorithm to the average relative error of a random set of 3D points, aligned to the ground truth data with the optimal similarity or affine transformation.

5.3 Discussion

Not surprisingly, our results (Section 5.2 in particular) show that affine shape can be recovered more reliably than depth. We expect this to be the case since the computation of affine shape does not require knowledge of the aspect-ratio of the camera, and since it does not require the computation of the square root of the Gramian matrix G .

Pt. #	Pose Z	Rigid Invar.		Aff. Invar.	
		Z	Err (%)	Z	Err (%)
1	14.4	16.8	16.3	14.4	0.2
2	15.1	15.1	-0.0	15.1	-0.1
3	14.5	16.3	12.5	14.3	-1.1
4	13.5	16.0	18.4	12.3	-9.1
5	21.7	23.7	9.6	21.8	0.9
6	18.8	20.1	7.0	18.4	-2.3
7	21.5	20.7	-4.0	22.0	2.3
8	20.0	23.7	18.0	19.8	-1.3
9	21.6	21.5	-0.5	22.3	2.9
10	21.0	22.5	7.3	21.8	4.1
11	21.6	20.1	-7.0	22.7	4.9
12	21.0	21.0	0.3	22.2	6.0
ave.			8.4%		2.9%

Table 2: The relative errors in depth computation using our invariant algorithm, for affine and rigid shape.

Rigid Invar.	Rigid random	Aff. Invar	Aff. random
8.4%	27.6%	2.9%	23.3%

Table 3: The mean relative errors in depth computation.

The sequence discussed in Section 5.1 was taken at a relatively large distance between the camera and the object (the depth values of the points varied from 550 to 700 mms). The weak perspective assumption therefore gave a good approximation. This sequence is typical of a recognition task. Under these conditions, which lend themselves favorably to the weak perspective approximation, our algorithm clearly performs very well. When compared with the other two algorithms, our algorithm is more efficient in its time complexity, it is simpler to implement, and it does not make any assumption on the type of motion (namely, it does not use the knowledge that the motion is rotational).

The sequence discussed in Section 5.2 had very large perspective distortions (the depth values of the points varied from 13 to 33 feet). Moreover, the sequence was obtained with a wide-lens camera, which lead to distortions in the image coordinates of points at the periphery. This sequence is more typical of a navigation task. Under these conditions, which do not lend themselves favorably to the weak perspective approximation, our algorithm is not accurate. The accuracy is sufficient for tasks which require only relative depth (e.g., obstacle avoidance), or less precise reconstruction of the environment. Note, however, that even algorithms which use the perspective projection model do not necessarily perform better with such sequences (compare with the results for a similar sequence reported in [10]).

In this last sequence, the computation of invariant shape using 8 frames or 16 frames lead to rather similar results for the affine shape matrix and the Gramian matrix. However, in the second case the computed Gramian matrix was not positive-definite, and therefore we could not compute depth. This demonstrates how the computation of depth is more sensitive to errors than the computation of the Euclidean representation. For the same reason, the affine reconstruction was an order of magnitude closer to the ground truth values than a set of random points, whereas

the depth reconstruction had an average error only 3 times smaller than a set of random points.

References

- [1] R. C. Bolles, H. H. Baker, and D. H. Marimont. Epipolar-plane image analysis: An approach to determining structure from motion. *International Journal of Computer Vision*, 1(1):7–55, 1987.
- [2] W. E. L. Grimson, D. P. Huttenlocher, and D. W. Jacobs. A study of affine matching with bounded sensor error. In G. Sandini, editor, *Computer Vision – ECCV92*, pages 291–306, Berlin, May 1992. Springer-Verlag.
- [3] B. K. P. Horn, H. M. Hilden, and S. Negahdaripour. Closed-form solution of absolute orientation using orthonormal matrices. *Journal of the Optical Society of America A*, 5(7):1127–1135, July 1988.
- [4] B. K. P. Horn. Relative orientation. *International Journal of Computer Vision*, 4(1):59–78, 1990.
- [5] D. J. Heeger and A. Jepson. Visual perception of three-dimensional motion. Technical Report 124, MIT Media Laboratory, Cambridge, Ma, December 1989.
- [6] R. Kumar and A. R. Hanson. Sensitivity of the pose refinement problem to accurate estimation of camera parameters. In *Proceedings of the 3rd International Conference on Computer Vision*, pages 365–369, Osaka, Japan, 1990. IEEE, Washington, DC.
- [7] J. J. Koenderink and A. J. van Doorn. Affine structure from motion. *Journal of the Optical Society of America*, 8(2):377–385, 1991.
- [8] D. T. Lawton. Processing translational motion sequences. *Computer Graphics and Image Processing*, 22:116–144, 1983.
- [9] H. C. Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293:133–135, September 1981.
- [10] H. S. Sawhney, J. Oliensis, and A. R. Hanson. Description and reconstruction from image trajectories of rotational motion. In *Proceedings of the 3rd International Conference on Computer Vision*, pages 494–498, Osaka, Japan, 1990. IEEE, Washington, DC.
- [11] C. Tomasi and T. Kanade. Shape and motion from image streams: a factorization method - 3. detection and tracking of point features. Technical Report CMU-CS-91-132, Carnegie Mellon University, Pittsburgh, PA, April 1991.
- [12] C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: a factorization method. *International Journal of Computer Vision*, 9(2):137–154, 1992.
- [13] S. Ullman and R. Basri. Recognition by linear combinations of models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13:992–1006, 1991.
- [14] S. Ullman. *The Interpretation of Visual Motion*. The MIT Press, Cambridge, MA, 1979.
- [15] D. Weinshall. Model-based invariants for 3D vision. *International Journal on Computer Vision*, 10(1):27–42,

1993.

- [16] D. Weinshall and C. Tomasi. Linear and incremental acquisition of invariant shape models from image sequences. RC 18549 (81133), IBM T. J. Watson Research Center, 1992.

Pt. #	Pose Z	Rigid Invar.		Aff. Invar.		Rot.		2-frm		2-frm, Ave.	
		Z	E (%)	Z	E (%)	Z	E (%)	Z	E (%)	Z	E (%)
1	591.4	587.8	-0.6	591.3	-0.0	588.9	-0.4	613.9	3.8	591.7	0.1
2	666.3	669.7	0.5	662.0	-0.7	665.8	-0.1	694.4	4.2	666.4	0.0
3	621.8	618.3	-0.6	621.9	0.0	617.8	-0.6	648.4	4.3	624.9	0.5
4	640.7	642.2	0.2	640.2	-0.1	635.0	-0.9	667.5	4.2	641.5	0.1
5	637.7	633.7	-0.6	637.0	-0.1	637.7	0.0	665.0	4.3	639.6	0.3
6	647.9	647.6	-0.1	647.0	-0.1	650.9	0.5	679.2	4.8	651.7	0.6
7	656.6	656.5	-0.0	654.3	-0.3	661.9	0.8	687.5	4.7	658.8	0.3
8	640.0	640.2	0.0	639.7	-0.0	653.8	2.2	668.0	4.4	642.3	0.4
9	709.7	708.8	-0.1	706.5	-0.5	700.7	-1.3	744.8	5.0	714.4	0.7
10	614.8	615.4	0.1	618.2	0.5	603.6	-1.8	644.1	4.8	618.5	0.6
11	602.3	602.6	0.0	601.5	-0.1	606.2	0.6	626.9	4.1	604.8	0.4
12	628.9	631.2	0.4	626.3	-0.4	636.5	1.2	655.3	4.2	630.5	0.2
ave.			0.27%		0.23%		0.86%		4.4%		0.35%

Table 1: Comparison of the relative errors in depth computation using our algorithm (rigid and affine shape separately), with two other algorithms. The average of the absolute value of the relative errors is listed at the bottom for each algorithm.