



A Closed-Form Solution to Non-Rigid Shape and Motion Recovery

JING XIAO

Epson Palo Alto Laboratory, 3145 Porter Drive, Palo Alto, CA 94304, USA

xiaoj@erd.epson.com,

jxiao@cs.cmu.edu (www.cs.cmu.edu/~jxiao)

JINXIANG CHAI AND TAKEO KANADE

Robotics Institute, CMU, 5000 Forbes Avenue, Pittsburgh, PA 15213, USA

jchai@cs.cmu.edu

tk@cs.cmu.edu

Received September 21, 2004; Revised April 19, 2005; Accepted May 3, 2005

First online version published in February, 2006

Abstract. Recovery of three dimensional (3D) shape and motion of non-static scenes from a monocular video sequence is important for applications like robot navigation and human computer interaction. If every point in the scene randomly moves, it is impossible to recover the non-rigid shapes. In practice, many non-rigid objects, e.g. the human face under various expressions, deform with certain structures. Their shapes can be regarded as a weighted combination of certain shape bases. Shape and motion recovery under such situations has attracted much interest. Previous work on this problem (Bregler, C., Hertzmann, A., and Biermann, H. 2000. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*; Brand, M. 2001. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*; Torresani, L., Yang, D., Alexander, G., and Bregler, C. 2001. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*) utilized only orthonormality constraints on the camera rotations (*rotation constraints*). This paper proves that using only the rotation constraints results in ambiguous and invalid solutions. The ambiguity arises from the fact that the shape bases are not unique. An arbitrary linear transformation of the bases produces another set of eligible bases. To eliminate the ambiguity, we propose a set of novel constraints, *basis constraints*, which uniquely determine the shape bases. We prove that, under the weak-perspective projection model, enforcing both the basis and the rotation constraints leads to a closed-form solution to the problem of non-rigid shape and motion recovery. The accuracy and robustness of our closed-form solution is evaluated quantitatively on synthetic data and qualitatively on real video sequences.

Keywords: non-rigid structure from motion, shape bases, rotation constraint, ambiguity, basis constraints, closed-form solution

1. Introduction

Many years of work in structure from motion have led to significant successes in recovery of 3D shapes and motion estimates from 2D monocular videos. Many reliable methods have been proposed for reconstruction of static scenes (Tomasi and Kanade, 1992; Zhang et al., 1995; Triggs, 1996; Hartley

and Zisserman, 2000; Mahamud and Hebert, 2000). However, most biological objects and natural scenes vary their shapes: expressive faces, people walking beside buildings, etc. Recovering the structure and motion of these non-rigid objects is a challenging task. The effects of rigid motion, i.e. 3D rotation and translation, and non-rigid shape deformation, e.g. stretching, are coupled together in the image

measurements. While it is impossible to reconstruct the shape if the scene deforms arbitrarily, in practice, many non-rigid objects, e.g. the human face under various expressions, deform with a class of structures.

One class of solutions model non-rigid object shapes as weighted combinations of certain shape bases that are pre-learned by off-line training (Bascle and Blake, 1998; Blanz and Vetter, 1999; Brand and Bhotika, 2001; Gokturk et al., 2001). For instance, the geometry of a face is represented as a weighted combination of shape bases that correspond to various facial deformations. Then the recovery of shape and motion is simply a model fitting problem. However, in many applications, e.g. reconstruction of a scene consisting of a moving car and a static building, the shape bases of the dynamic structure are difficult to obtain before reconstruction.

Several approaches have been proposed to solve the problem without a prior model (Bregler et al., 2000; Torresani et al., 2001; Brand, 2001). Instead, they treat the model, i.e. shape bases, as part of the unknowns to be computed. They try to recover not only the non-rigid shape and motion, but also the shape model. This class of approaches so far has utilized only the orthonormality constraints on camera rotations (*rotation constraints*) to solve the problem. However, as shown in this paper, enforcing only the rotation constraints leads to ambiguous and invalid solutions. These approaches thus cannot guarantee the desired solution. They have to either require a priori knowledge on shape and motion, e.g. constant speed (Han and Kanade, 2001), or need non-linear optimization that involves large number of variables and hence requires a good initial estimate (Torresani et al., 2001; Brand, 2001).

Intuitively, the above ambiguity arises from the non-uniqueness of the shape bases: an arbitrary linear transformation of a set of shape bases yields a new set of eligible bases. Once the bases are determined uniquely, the ambiguity is eliminated. Therefore, instead of imposing only the rotation constraints, we identify and introduce another set of constraints on the shape bases (*basis constraints*), which implicitly determine the bases uniquely. This paper proves that, under the weak-perspective projection model, when both the basis and rotation constraints are imposed, a linear closed-form solution to the problem of non-rigid shape and motion recovery is achieved. Accordingly we develop a factorization method that applies both the metric constraints to compute the closed-form solution for the non-rigid shape, motion, and shape bases.

2. Previous Work

Recovering 3D object structure and motion from 2D image sequences has a rich history. Various approaches have been proposed for different applications. The discussion in this section will focus on the factorization techniques, which are most closely related to our work.

The factorization method was originally proposed by Tomasi and Kanade (1992). First it applies the rank constraint to factorize a set of feature locations tracked across the entire sequence. Then it uses the orthonormality constraints on the rotation matrices to recover the scene structure and camera rotations in one step. This approach works under the orthographic projection model. Poelman and Kanade (1997) extended it to work under the weak perspective and para-perspective projection models. Triggs (1996) generalized the factorization method to the recovery of scene geometry and camera motion under the perspective projection model. These methods work for static scenes.

The factorization technique has been extended to the non-rigid cases (Costeira and Kanade, 1998; Han and Kanade, 2000; Wolf and Shashua, 2001; Kanatani, 2001; Wolf and Shashua, 2002; Vidal and Soatto, 2002; Zelnik-Manor and Irani, 2003; Sugaya and Kanatani, 2004; Vidal and Hartley, 2004). Costeira and Kanade (1998) proposed a method to reconstruct the structure of multiple independently moving objects under the weak-perspective camera model. This method first separates the objects according to their respective factorization characteristics and then individually recovers their shapes using the original factorization technique. The similar reconstruction problem was also studied in Kanatani (2001), Zelnik-Manor and Irani (2003), and Sugaya and Kanatani (2004). Wolf and Shashua (2001) derived a geometrical constraint, called the segmentation matrix, to reconstruct a scene containing two independently moving objects from two perspective views. Vidal and his colleagues (Vidal et al., 2002; Vidal and Hartley, 2004) generalized this approach for reconstructing dynamic scenes containing multiple independently moving objects. In cases where the dynamic scenes consist of both static objects and objects moving straight along respective directions, Han and Kanade (2000) proposed a weak-perspective reconstruction method that achieves a closed-form solution, under the assumption of constant moving velocities. A more generalized solution to reconstructing the shapes that deform at constant velocity is presented in Wolf and Shashua (2002).

Bregler and his colleagues (Bregler et al., 2000) for the first time incorporated the basis representation of non-rigid shapes into the reconstruction process. By analyzing the low rank of the image measurements, they proposed a factorization-based method that enforces the orthonormality constraints on camera rotations to reconstruct the non-rigid shape and motion. Torresani and his colleagues (Torresani et al., 2001) extended the method in Bregler et al. (2000) to a trilinear optimization approach. At each step, two of the three types of unknowns, bases, coefficients, and rotations, are fixed and the remaining one is updated. The method in Bregler et al. (2000) is used to initialize the optimization process. Brand (2001) proposed a similar non-linear optimization method that uses an extension of the method in Bregler et al. (2000) for initialization. All the three methods enforce the rotation constraints alone and thus cannot guarantee the optimal solution. Note that because both of the non-linear optimization methods involve a large number of variables, e.g. the number of coefficients equals the product of the number of the images and the number of the shape bases, their performances greatly rely on the quality of the initial estimates.

3. Problem Statement

Given 2D locations of P feature points across F frames, $\{(v, v)_{fp}^T | f = 1, \dots, F, p = 1, \dots, P\}$, our goal is to recover the motion of the non-rigid object relative to the camera, including rotations $\{R_f | f = 1, \dots, F\}$ and translations $\{\mathbf{t}_f | f = 1, \dots, F\}$, and its 3D deforming shapes $\{(x, y, z)_{fp}^T | f = 1, \dots, F, p = 1, \dots, P\}$. Throughout this paper, we assume:

- the deforming shapes can be represented as weighted combinations of shape bases;
- the 3D structure and the camera motion are non-degenerate;
- the camera projection model is the weak-perspective projection model.

We follow the representation of Blanz and Vetter (1999) and Bregler et al. (2000). The non-rigid shapes are represented as weighted combinations of K shape bases $\{B_i, i = 1, \dots, K\}$. The bases are $3 \times P$ matrices controlling the deformation of P points. Then the 3D coordinate of the point p at the frame f is

$$\mathbf{X}_{fp} = (x, y, z)_{fp}^T = \sum_{i=1}^K c_{fi} \mathbf{b}_{ip} \quad f = 1, \dots, F, p = 1, \dots, P \quad (1)$$

where \mathbf{b}_{ip} is the p th column of B_i and c_{fi} is its combination coefficient at the frame f . The image coordinate of \mathbf{X}_{fp} under the weak perspective projection model is

$$\mathbf{x}_{fp} = (u, v)_{fp}^T = s_f(R_f \cdot \mathbf{X}_{fp} + \mathbf{t}_f) \quad (2)$$

where R_f stands for the first two rows of the f th camera rotation and $\mathbf{t}_f = [t_{fx} t_{fy}]^T$ is its translation relative to the world origin. s_f is the scalar of the weak perspective projection.

Replacing \mathbf{X}_{fp} using Eq. (1) and absorbing s_f into c_{fi} and \mathbf{t}_f , we have

$$\mathbf{x}_{fp} = (c_{f1} R_f \dots c_{fK} R_f) \cdot \begin{pmatrix} \mathbf{b}_{1p} \\ \vdots \\ \mathbf{b}_{Kp} \end{pmatrix} + \mathbf{t}_f \quad (3)$$

Suppose the image coordinates of all P feature points across F frames are known. We form a $2F \times P$ measurement matrix W by stacking all image coordinates. Then $W = MB + T[11 \dots 1]$, where M is a $2F \times 3K$ scaled rotation matrix, B is a $3K \times P$ bases matrix, and T is a $2F \times 1$ translation vector,

$$M = \begin{pmatrix} c_{11} R_1 & \dots & c_{1K} R_1 \\ \vdots & \vdots & \vdots \\ c_{F1} R_F & \dots & c_{FK} R_F \end{pmatrix}, \quad B = \begin{pmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \vdots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \end{pmatrix}, \quad T = (\mathbf{t}_1^T \dots \mathbf{t}_F^T)^T \quad (4)$$

As in Han and Kanade (2000) and Bregler et al. (2000), since all points on the shape share the same translation, we position the world origin at the scene center and compute the translation vector by averaging the image projections of all points. This step does not introduce ambiguities, provided that the correspondences of all the points across all the images are given, i.e. there are no missing data. We then subtract it from W and obtain the *registered* measurement matrix $\tilde{W} = MB$.

Under the non-degenerate cases, the $2F \times 3K$ scaled rotation matrix M and the $3K \times P$ shape bases matrix B are both of full rank, respectively $\min\{2F, 3K\}$ and $\min\{3K, P\}$. Their product, \tilde{W} , is of rank $\min\{3K, 2F, P\}$. In practice, the frame number F and point number P are usually much larger than

the basis number K such that $2F > 3K$ and $P > 3K$. Thus the rank of \tilde{W} is $3K$ and K is determined by $K = \frac{\text{rank}(\tilde{W})}{3}$. We then factorize \tilde{W} into the product of a $2F \times 3K$ matrix \tilde{M} and a $3K \times P$ matrix \tilde{B} , using Singular Value Decomposition (SVD). This decomposition is only determined up to a non-singular $3K \times 3K$ linear transformation. The true scaled rotation matrix M and bases matrix B are of the form,

$$M = \tilde{M} \cdot G, \quad B = G^{-1} \cdot \tilde{B} \quad (5)$$

where the non-singular $3K \times 3K$ matrix G is called the *corrective transformation* matrix. Once G is determined, M and B are obtained and thus the rotations, shape bases, and combination coefficients are recovered.

All the procedures above, except obtaining G , are standard and well-understood (Blaž and Vetter, 1999; Bregler et al., 2000). The problem of nonrigid shape and motion recovery is now reduced to: given the measurement matrix W , how can we compute the *corrective transformation* matrix G ?

4. Metric Constraints

G is made up of K triple-columns, denoted as $g_k, k = 1, \dots, K$. Each of them is a $3K \times 3$ matrix. They are independent on each other because G is non-singular. According to Eqs. (4, 5), g_k satisfies,

$$\tilde{M} g_k = \begin{pmatrix} c_{1k} R_1 \\ \vdots \\ c_{Fk} R_F \end{pmatrix} \quad (6)$$

Then,

$$\begin{aligned} \tilde{M} g_k g_k^T \tilde{M}^T &= \begin{pmatrix} c_{1k}^2 R_1 R_1^T & c_{1k} c_{2k} R_1 R_2^T & \dots & c_{1k} c_{Fk} R_1 R_F^T \\ c_{1k} c_{2k} R_2 R_1^T & c_{2k}^2 R_2 R_2^T & \dots & c_{2k} c_{Fk} R_2 R_F^T \\ \vdots & \vdots & \ddots & \vdots \\ c_{1k} c_{Fk} R_F R_1^T & c_{2k} c_{Fk} R_F R_2^T & \dots & c_{Fk}^2 R_F R_F^T \end{pmatrix} \end{aligned} \quad (7)$$

We denote $g_k g_k^T$ by Q_k , a $3K \times 3K$ symmetric matrix. Once Q_k is determined, g_k can be computed using SVD. To compute Q_k , two types of metric constraints are available and should be imposed: *rotation constraints*

and *basis constraints*. While using only the rotation constraints (Bregler et al., 2000; Brand, 2001) leads to ambiguous and invalid solutions, enforcing both sets of constraints results in a linear closed-form solution for Q_k .

4.1. Rotation Constraints

The orthonormality constraints on the rotation matrices are one of the most powerful metric constraints and they have been used in reconstructing the shape and motion for static objects (Tomasi and Kanade, 1992; Poelman and Kanade, 1997), multiple moving objects (Costeira and Kanade, 1998; Han and Kanade, 2000), and non-rigid deforming objects (Bregler et al., 2000; Torresani et al., 2001; Brand, 2001).

According to Eq. (7), we have,

$$\tilde{M}_{2i-1:2i} Q_k \tilde{M}_{2j-1:2j}^T = c_{ik} c_{jk} R_i R_j^T, \quad i, j = 1, \dots, F \quad (8)$$

where $\tilde{M}_{2i-1:2i}$ represents the i th bi-row of \tilde{M} . Due to orthonormality of the rotation matrices, we have,

$$\tilde{M}_{2i-1:2i} Q_k \tilde{M}_{2i-1:2i}^T = c_{ik}^2 \mathbf{I}_{2 \times 2}, \quad i = 1, \dots, F \quad (9)$$

where $\mathbf{I}_{2 \times 2}$ is a 2×2 identity matrix. The two diagonal elements of Eq. (9) yield one linear constraints on Q_k , since c_{ik} is unknown. The two off-diagonal constraints are identical, because Q_k is symmetric. For all F frames, we obtain $2F$ linear constraints as follows,

$$\tilde{M}_{2i-1} Q_k \tilde{M}_{2i-1}^T - \tilde{M}_{2i} Q_k \tilde{M}_{2i}^T = 0, \quad i = 1, \dots, F \quad (10)$$

$$\tilde{M}_{2i-1} Q_k \tilde{M}_{2i}^T = 0, \quad i = 1, \dots, F \quad (11)$$

We consider the entries of Q_k as variables and neglect the nonlinear constraints implicitly embedded in the formation of $Q_k = g_k g_k^T$, where g_k has only $9K$ free entries. Since Q_k is symmetric, it contains $\frac{(9K^2+3K)}{2}$ independent variables. It appears that, when enough images are given, i.e. $2F \geq \frac{(9K^2+3K)}{2}$, the rotation constraints in Eqs. (10, 11) should be sufficient to determine Q_k via the linear least-square method. However, it is not true in general. We will show that most of the rotation constraints are redundant and they are inherently insufficient to resolve Q_k .

4.2. Why are Rotation Constraints Not Sufficient?

Under specific assumptions like static scene or constant speed of deformation, the orthonormality constraints are sufficient to reconstruct the 3D shapes and camera rotations (Tomasi and Kanade, 1992; Han and Kanade, 2000). In general cases, however, no matter how many images or feature points are given, the solution of the rotation constraints in Eqs. (10, 11) is inherently ambiguous.

Definition 1. A $3K \times 3K$ symmetric matrix Y is called a block-skew-symmetric matrix, if all the diagonal 3×3 blocks are zero matrices and each off-diagonal 3×3 block is a skew symmetric matrix.

$$Y_{ij} = \begin{pmatrix} 0 & y_{ij1} & y_{ij2} \\ -y_{ij1} & 0 & y_{ij3} \\ -y_{ij2} & -y_{ij3} & 0 \end{pmatrix} = -Y_{ij}^T = Y_{ji}^T, \quad i \neq j \quad (12)$$

$$Y_{ii} = 0_{3 \times 3}, \quad i, j = 1, \dots, K \quad (13)$$

Each off-diagonal block consists of 3 independent elements. Because Y is skew-symmetric and has $\frac{K(K-1)}{2}$ independent off-diagonal blocks, it includes $\frac{3K(K-1)}{2}$ independent elements.

Definition 2. A $3K \times 3K$ symmetric matrix Z is called a block-scaled-identity matrix, if each 3×3 block is a scaled identity matrix, i.e. $Z_{ij} = \lambda_{ij} \mathbf{I}_{3 \times 3}$, where λ_{ij} is the only variable.

Because Z is symmetric, the total number of variables in Z equals the number of independent blocks, $\frac{K(K+1)}{2}$.

Theorem 1. The general solution of the rotation constraints in Eqs. (10, 11) is GH_kG^T , where G is the desired corrective transformation matrix, and H_k is the summation of an arbitrary block-skew-symmetric matrix and an arbitrary block-scaled-identity matrix, given a minimum of $\frac{K^2+K}{2}$ images containing different (non-proportional) shapes and another $\frac{K^2+K}{2}$ images with non-degenerate (non-planar) rotations.

Proof: Let us denote \tilde{Q}_k as the general solution of Eqs. (10, 11). Since G is a non-singular square matrix, \tilde{Q}_k can be represented as GH_kG^T , where $H_k = G^{-1}\tilde{Q}_kG^{-T}$. We then prove that H_k must be

the summation of a block-skew-symmetric matrix and a block-scaled-identity matrix.

According to Eqs. (5, 9),

$$\begin{aligned} c_{ik}^2 \mathbf{I}_{2 \times 2} &= \tilde{M}_{2i-1:2i} \tilde{Q}_k \tilde{M}_{2i-1:2i}^T \\ &= \tilde{M}_{2i-1:2i} G H_k G^T \tilde{M}_{2i-1:2i}^T \\ &= M_{2i-1:2i} H_k M_{2i-1:2i}^T, \quad i = 1, \dots, F \end{aligned} \quad (14)$$

H_k consists of $K^2 3 \times 3$ blocks, denoted as H_{kmn} , $m, n = 1, \dots, K$. Combining Eqs. (4) and (14), we have,

$$\begin{aligned} R_i \sum_{m=1}^K (c_{im}^2 H_{kmm} + \sum_{n=m+1}^K c_{im} c_{in} (H_{kmn} + H_{knn}^T)) R_i^T \\ = c_{ik}^2 \mathbf{I}_{2 \times 2}, \quad i = 1, \dots, F \end{aligned} \quad (15)$$

Denote the 3×3 symmetric matrix $\sum_{m=1}^K (c_{im}^2 H_{kmm} + \sum_{n=m+1}^K c_{im} c_{in} (H_{kmn} + H_{knn}^T))$ by Γ_{ik} . Then Eq. (15) becomes $R_i \Gamma_{ik} R_i^T = c_{ik}^2 \mathbf{I}_{2 \times 2}$ be its homogeneous solution, i.e. $R_i \tilde{\Gamma}_{ik} R_i^T = \mathbf{0}_{2 \times 2}$. So we are looking for a symmetric matrix $\tilde{\Gamma}_{ik}$ in the null space of the map $X \mapsto R_i X R_i^T$. Because the two rows of the 2×3 matrix R_i are orthonormal, we have,

$$\tilde{\Gamma}_{ik} = r_{i3}^T \delta_{ik} + \delta_{ik}^T r_{i3} \quad (16)$$

where r_{i3} is a unitary 1×3 vector that are orthogonal to both rows of R_i . δ_{ik} is an arbitrary 1×3 vector, indicating the null space of the map $X \mapsto R_i X R_i^T$ has 3 degrees of freedom. $\Gamma_{ik} = c_{ik}^2 \mathbf{I}_{3 \times 3}$ is a particular solution of $R_i \Gamma_{ik} R_i^T = c_{ik}^2 \mathbf{I}_{2 \times 2}$. Thus the general solution of Eq. (15) is,

$$\begin{aligned} \sum_{m=1}^K (c_{im}^2 H_{kmm} + \sum_{n=m+1}^K c_{im} c_{in} (H_{kmn} + H_{knn}^T)) \\ = \Gamma_{ik} = c_{ik}^2 \mathbf{I}_{3 \times 3} + \alpha_{ik} \tilde{\Gamma}_{ik}, \quad i = 1, \dots, F \end{aligned} \quad (17)$$

where α_{ik} is an arbitrary scalar.

Let us denote $N = \frac{K^2+K}{2}$. Each image yields a set of N coefficients $(c_{i1}^2, c_{i1}c_{i2}, \dots, c_{i1}c_{iK}, c_{i2}^2, c_{i2}c_{i3}, \dots, c_{iK-1}c_{iK}, c_{iK}^2)$ in Eq. (17). We call them the quadratic-form coefficients. The quadratic-form coefficients associated with images with non-proportional shapes are generally independent. In our test, the samples of randomly generated N different (non-proportional) shapes always yield independent N sets of quadratic-form coefficients. Thus given N images with different (non-proportional) shapes, the associated N quadratic-form

coefficient sets span the space of the quadratic-form coefficient in Eq. (17). The linear combinations of these N sets compose the quadratic-form coefficients associated with any other images. Since H_{kmm} and H_{kmn} are common in all images, for any additional image $j = 1, \dots, F$, Γ_{jk} is described as follows,

$$\begin{aligned}\Gamma_{jk} &= \sum_{l=1}^N \lambda_{lk} \Gamma_{lk} \\ \Rightarrow c_{jk}^2 \mathbf{I}_{3 \times 3} + \alpha_{jk} \tilde{\Gamma}_{jk} &= \sum_{l=1}^N \lambda_{lk} (c_{lk}^2 \mathbf{I}_{3 \times 3} + \alpha_{lk} \tilde{\Gamma}_{lk}) \\ \Rightarrow \alpha_{jk} \tilde{\Gamma}_{jk} &= \sum_{l=1}^N \lambda_{lk} (\alpha_{lk} \tilde{\Gamma}_{lk})\end{aligned}\quad (18)$$

where λ_{lk} is the combination weights. We substitute $\tilde{\Gamma}_{jk}$ and $\tilde{\Gamma}_{lk}$ by Eq. (16) and absorb α_{jk} and α_{lk} into δ_{jk} and δ_{lk} respectively due to their arbitrary values. Eq. (18) is then rewritten as,

$$r_{j3}^T \delta_{jk} + \delta_{jk}^T r_{j3} - \sum_{l=1}^N \lambda_{lk} (r_{l3}^T \delta_{lk} + \delta_{lk}^T r_{l3}) = 0 \quad (19)$$

Due to symmetry, Eq. (19) yields 6 linear constraints on δ_{lk} , $l = 1, \dots, N$ and δ_{jk} that in total include $3N + 3$ variables. These constraints depend on the rotations in the images. Given x additional images with non-degenerate (non-planar) rotations, we obtain $6x$ different linear constraints and $3N + 3x$ free variables. When $6x \geq (3N + 3x)$, i.e. $x \geq N$, these constraints lead to a unique solution, $\delta_{lk} = \delta_{jk} = 0$, $l = 1, \dots, N$, $j = 1, \dots, x$. Therefore, given a minimum of N images with non-degenerate rotations together with the N images containing different shapes, Eq. (17) can be rewritten as follows,

$$\begin{aligned}\Gamma_{ik} &= \sum_{m=1}^K (c_{im}^2 H_{kmm} + \sum_{n=m+1}^K c_{im} c_{in} (H_{kmn} + H_{kmn}^T)) \\ &= c_{ik}^2 \mathbf{I}_{3 \times 3}, \quad i = 1, \dots, 2N\end{aligned}\quad (20)$$

Eq. (20) completely depend on the coefficients. According to Eq. (18, 20), for any image $j = 2N + 1, \dots, F$, $\Gamma_{jk} = \sum_{l=1}^N \lambda_{lk} \Gamma_{lk} = \sum_{l=1}^N \lambda_{lk} c_{lk}^2 \mathbf{I}_{3 \times 3} = c_{jk}^2 \mathbf{I}_{3 \times 3}$. Therefore Eq. (20) is satisfied for all F images. Denote $(H_{kmn} + H_{kmn}^T)$ by Θ_{kmn} . Because the right side of Eq. (20) is a scaled identity matrix, for each off-diagonal element h_{kmn}^o and θ_{kmn}^o , we achieve the following linear equation,

$$\sum_{m=1}^K (c_{im}^2 h_{kmn}^o + \sum_{n=m+1}^K c_{im} c_{in} \theta_{kmn}^o) = 0, \quad i = 1, \dots, F \quad (21)$$

Using the N given images containing different shapes, Eq. (21) leads to a non-singular linear equation set on h_{kmn}^o and θ_{kmn}^o . The right sides of the equations are

zeros. Thus the solution is a zero vector, i.e. the off-diagonal elements of H_{kmm} and Θ_{kmn} are all zeros. Similarly, we can derive the constraint as Eq. (21) on the difference between any two diagonal elements. Therefore the difference is zero, i.e. the diagonal elements are all equivalent. We thus have,

$$H_{kmm} = \lambda_{kmm} \mathbf{I}_{3 \times 3}, \quad m = 1, \dots, K \quad (22)$$

$$\begin{aligned}H_{kmn} + H_{kmn}^T &= \lambda_{kmn} \mathbf{I}_{3 \times 3}, \quad m = 1, \dots, K, \\ n &= m + 1, \dots, K\end{aligned}\quad (23)$$

where λ_{kmm} and λ_{kmn} are arbitrary scalars. According to Eq. (22), the diagonal block H_{kmm} is a scaled identity matrix. Due to Eq. (23), $H_{kmn} - \frac{\lambda_{kmn}}{2} \mathbf{I}_{3 \times 3} = -(H_{kmn} - \frac{\lambda_{kmn}}{2} \mathbf{I}_{3 \times 3})^T$, i.e. $H_{kmn} - \frac{\lambda_{kmn}}{2} \mathbf{I}_{3 \times 3}$ is skew-symmetric. Thus the off-diagonal block H_{kmn} equals the summation of a scaled identity block, $\frac{\lambda_{kmn}}{2} \mathbf{I}_{3 \times 3}$, and a skew-symmetric block, $H_{kmn} - \frac{\lambda_{kmn}}{2} \mathbf{I}_{3 \times 3}$. Since λ_{kmm} and λ_{kmn} are arbitrary, the entire matrix H_k is the summation of an arbitrary block-skew-symmetric matrix and an arbitrary block-scaled-identity matrix, given a minimum of N images with non-degenerate rotations together with N images containing different shapes, in total $2N = K^2 + K$ images. \square

Let Y_k denote the block-skew-symmetric matrix and Z_k denote the block-scaled-identity matrix in H_k . Since Y_k and Z_k respectively contain $\frac{3K(K-1)}{2}$ and $\frac{K(K+1)}{2}$ independent elements, H_k include $2K^2 - K$ free elements, i.e. the solution of the rotation constraints has $2K^2 - K$ degrees of freedom. In rigid cases, i.e. $K = 1$, the degree of freedom is 1 ($2 \cdot 1^2 - 1 = 1$), meaning that the solution is unique up to a scale, as suggested in (Tomasi and Kanade, 1992).

For non-rigid objects, i.e. $K > 1$, the rotation constraints result in an ambiguous solution space. For the many combinations of Z_k and Y_k of which H_k are positive semi-definite, $\tilde{g}_k = G\sqrt{H_k}$ lead to the ambiguous solutions. Of course the space contains invalid solutions as well. Specifically, for those combinations of which H_k is not positive semi-definite, $\sqrt{H_k}$ is indefinite and $\tilde{g}_k = G\sqrt{H_k}$ refer to the invalid solutions. Such combinations indeed exist in the space. For example, when the block-scaled-identity matrix Z_k is zero, H_k equals the block-skew-symmetric matrix Y_k , which is not positive semi-definite.

4.3. Basis Constraints

The only difference between non-rigid and rigid situations is that the non-rigid shape is a weighted combination of certain shape bases. The rotation constraints are sufficient for recovering the rigid shapes, but they cannot determine a unique set of shape bases in the non-rigid cases. Instead any non-singular linear transformation applied on the bases leads to another set of eligible bases. Intuitively, the non-uniqueness of the bases results in the ambiguity of the solution by enforcing the rotation constraints alone. We thus introduce the basis constraints that determine a unique basis set and resolve the ambiguity.

Because the deformable shapes lie in a K -basis linear space, any K independent shapes in the space form an eligible bases set, i.e. an arbitrary shape in the space can be uniquely represented as a linear combination of these K independent shapes. Note that these independent bases are not necessarily orthogonal. We measure the independence of the K shapes using the condition number of the matrix that stacks the 3D coordinates of the shapes like the measurement matrix W . The condition number is infinitely big, if any of the K shapes is dependent on the other shapes, i.e. it can be described as a linear combination of the other shapes. If the K shapes are orthonormal, i.e. completely independent and equally dominant, the condition number achieves its minimum, 1. Between 1 and infinity, a smaller condition number means that the shapes are more independent and more equally dominant. The registered image measurements of these K shapes in \tilde{W} are,

$$\begin{pmatrix} \tilde{W}_1 \\ \tilde{W}_2 \\ \vdots \\ \tilde{W}_K \end{pmatrix} = \begin{pmatrix} R_1 & 0 & \dots & 0 \\ 0 & R_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & R_K \end{pmatrix} \begin{pmatrix} S_1 \\ S_2 \\ \vdots \\ S_K \end{pmatrix} \quad (24)$$

where \tilde{W}_i and $S_i, i = 1, \dots, K$ are respectively the image measurements and the 3D coordinates of the K shapes. On the right side of Eq. (24), since the matrix composed of the rotations is orthonormal, the projecting process does not change the singular values of the 3D coordinate matrix. Therefore, the condition number of the measurement matrix has the same power as that of the 3D coordinate matrix to represent the independence of the K shapes. Accordingly we compute the condition number of the measurement matrix in each group of K images and identify a group of K images

with the smallest condition number. The 3D shapes in these K images are specified as the shape bases. Note that so far we have not recovered the bases explicitly, but decided in which frames they are located. This information implicitly determines a unique set of bases.

We denote the selected frames as the first K images in the sequence. The corresponding coefficients are,

$$\begin{aligned} c_{ii} &= 1, \quad i = 1, \dots, K \\ c_{ij} &= 0, \quad i, j = 1, \dots, K, \quad i \neq j \end{aligned} \quad (25)$$

Let Ω denote the set, $\{(i, j) | i = 1, \dots, K; i \neq j; j = 1, \dots, F\}$. According to Eqs. (8, 25), we have,

$$\tilde{M}_{2i-1} Q_k \tilde{M}_{2j-1}^T = \begin{cases} 1, & i = j = k \\ 0, & (i, j) \in \Omega \end{cases} \quad (26)$$

$$\tilde{M}_{2i} Q_k \tilde{M}_{2j}^T = \begin{cases} 1, & i = j = k \\ 0, & (i, j) \in \Omega \end{cases} \quad (27)$$

$$\tilde{M}_{2i-1} Q_k \tilde{M}_{2j}^T = 0, \quad (i, j) \in \Omega \text{ or } i = j = k \quad (28)$$

$$\tilde{M}_{2i} Q_k \tilde{M}_{2j-1}^T = 0, \quad (i, j) \in \Omega \text{ or } i = j = k \quad (29)$$

These $4F(K-1)$ linear constraints are called the basis constraints.

5. A Closed-Form Solution

Due to Theorem 1, enforcing the rotation constraints on Q_k leads to the ambiguous solution $GH_k G^T \cdot H_k$ consists of $K^2 3 \times 3$ blocks, $H_{kmn}, m, n = 1, \dots, K$. H_{kmn} contains four independent entries as follows,

$$H_{kmn} = \begin{pmatrix} h_1 & h_2 & h_3 \\ -h_2 & h_1 & h_4 \\ -h_3 & -h_4 & h_1 \end{pmatrix} \quad (30)$$

Lemma 1. H_{kmn} is a zero matrix if,

$$R_i H_{kmn} R_j^T = \mathbf{0}_{2 \times 2} \quad (31)$$

where R_i and R_j are non-degenerate 2×3 rotation matrices.

Proof: First we prove that the rank of H_{kmn} is at most 2. Due to the orthonormality of rotation matrices, from

Eq. (31), we have,

$$H_{kmn} = r_{i3}^T \delta_{ik} + \delta_{jk}^T r_{j3} = \begin{pmatrix} r_{i3}^T & \delta_{jk}^T \end{pmatrix} \begin{pmatrix} \delta_{ik} \\ r_{j3} \end{pmatrix} \quad (32)$$

where r_{i3} and r_{j3} respectively are the cross products of the two rows of R_i and those of R_j . δ_{ik} and δ_{jk} are arbitrary 1×3 vectors. Because both matrices on the right side of Eq. (32) are at most of rank 2, the rank of H_{kmn} is at most 2.

Next, we prove $h_1 = 0$. Since H_{kmn} is 3×3 matrix of rank 2, its determinant, $h_1(\sum_{i=1}^4 h_i^2)$, equals 0. Therefore $h_1 = 0$, i.e. H_{kmn} is a skew-symmetric matrix.

Finally we prove $h_2 = h_3 = h_4 = 0$. Denote the rows of R_i and R_j as r_{i1} , r_{i2} , r_{j1} , and r_{j2} respectively. Since $h_1 = 0$, we can rewrite Eq. (31) as follows,

$$\begin{pmatrix} r_{i1} \cdot (\mathbf{h} \times r_{j1}) & r_{i1} \cdot (\mathbf{h} \times r_{j2}) \\ r_{i2} \cdot (\mathbf{h} \times r_{j1}) & r_{i2} \cdot (\mathbf{h} \times r_{j2}) \end{pmatrix} = \mathbf{0}_{2 \times 2} \quad (33)$$

where $\mathbf{h} = (-h_4, h_3 - h_2)$. Equation (33) means that the vector \mathbf{h} is located in the intersection of the four planes determined by (r_{i1}, r_{j1}) , (r_{i1}, r_{j2}) , (r_{i2}, r_{j1}) , and (r_{i2}, r_{j2}) . Since R_i and R_j are non-degenerate rotation matrices, r_{i1} , r_{i2} , r_{j1} , and r_{j2} do not lie in the same plane, hence the four planes intersect at the origin, i.e. $\mathbf{h} = (-h_4, h_3 - h_2) = \mathbf{0}_{1 \times 3}$. Therefore H_{kmn} is a zero matrix. \square

Using Lemma 1, we derive the following theorem,

Theorem 2. *Enforcing both the basis constraints and the rotation constraints results in a closed-form solution of Q_k , given a minimum of $\frac{K^2+K}{2}$ images containing different (non-proportional) shapes together with another $\frac{K^2+K}{2}$ images with non-degenerate (non-planar) rotations.*

Proof: According to Theorem 1, given a minimum of $\frac{K^2+K}{2}$ images containing different shapes and another $\frac{K^2+K}{2}$ images with non-degenerate rotations, using the rotation constraints we achieve the ambiguous solution of $Q_k = GH_k G^T$, where G is the desired corrective transformation matrix and H_k is the summation of an arbitrary block-skew-symmetric matrix and an arbitrary block-scaled-identity matrix.

Due to the basis constraints, replacing Q_k in Eq. (26–29) with $GH_k G^T$,

$$\tilde{M}_{2k-1:2k} GH_k G^T \tilde{M}_{2k-1:2k}^T$$

$$= M_{2k-1:2k} H_k M_{2k-1:2k}^T = \mathbf{I}_{2 \times 2} \quad (34)$$

$$\begin{aligned} \tilde{M}_{2i-1:2i} G H_k G^T \tilde{M}_{2j-1:2j}^T &= M_{2i-1:2i} H_k M_{2j-1:2j}^T \\ &= \mathbf{0}_{2 \times 2}, \quad i = 1, \dots, K, \quad j = 1, \dots, F, \quad i \neq k \end{aligned} \quad (35)$$

From Eq. (4), we have,

$$M_{2i-1:2i} H_k M_{2j-1:2j}^T = \sum_{m=1}^K \sum_{n=1}^K c_{im} c_{jn} R_i H_{kmn} R_j^T, \quad i, j = 1, \dots, F \quad (36)$$

where H_{kmn} is the 3×3 block of H_k . According to Eq. (25),

$$M_{2i-1:2i} H_k M_{2j-1:2j}^T = R_i H_{kij} R_j^T, \quad i, j = 1, \dots, K \quad (37)$$

Combining Eqs. (34, 35) and (37), we have,

$$R_k H_{kkk} R_k^T = \mathbf{I}_{2 \times 2} \quad (38)$$

$$R_i H_{kij} R_j^T = \mathbf{0}_{2 \times 2}, \quad i, j = 1, \dots, K, \quad i \neq k \quad (39)$$

According to Eq. (22), the k_{th} diagonal block of H_k , H_{kkk} , equals $\lambda_{kkk} \mathbf{I}_{3 \times 3}$. Thus by Eq. (38), $\lambda_{kkk} = 1$ and $H_{kkk} = \mathbf{I}_{3 \times 3}$. We denote K of the $\frac{K^2+K}{2}$ given images with non-degenerate rotations as the first K images in the sequence. Due to Lemma 1 and Eq. (39), H_{kij} , $i, j = 1, \dots, K, i \neq k$, are zero matrices. Since H_k is symmetric, the other blocks in H_k , H_{kij} , $i = k, j = 1, \dots, K, j \neq k$, are also zero matrices. Thus,

$$\begin{aligned} GH_k G^T &= (g_1 \dots g_K) H_k (g_1 \dots g_K)^T \\ &= (0 \dots g_k \dots 0) (g_1 \dots g_K)^T \\ &= g_k g_k^T = Q_k \end{aligned} \quad (40)$$

i.e. a closed-form solution of the desired Q_k has been achieved. \square

According to Theorem 2, we compute $Q_k = g_k g_k^T$, $k = 1, \dots, K$, by solving the linear equations, Eqs. (10–11, 26–29), via the least square methods. We then recover g_k by decomposing Q_k via SVD. The decomposition of Q_k is up to an arbitrary 3×3 orthonormal transformation Φ_k , since $(g_k \Phi_k)(g_k \Phi_k)^T$ also equals Q_k . This ambiguity arises from the fact that g_1, \dots, g_K are estimated independently under different coordinate systems. To resolve the ambiguity, we need to transform g_1, \dots, g_K to be under a common reference coordinate system.

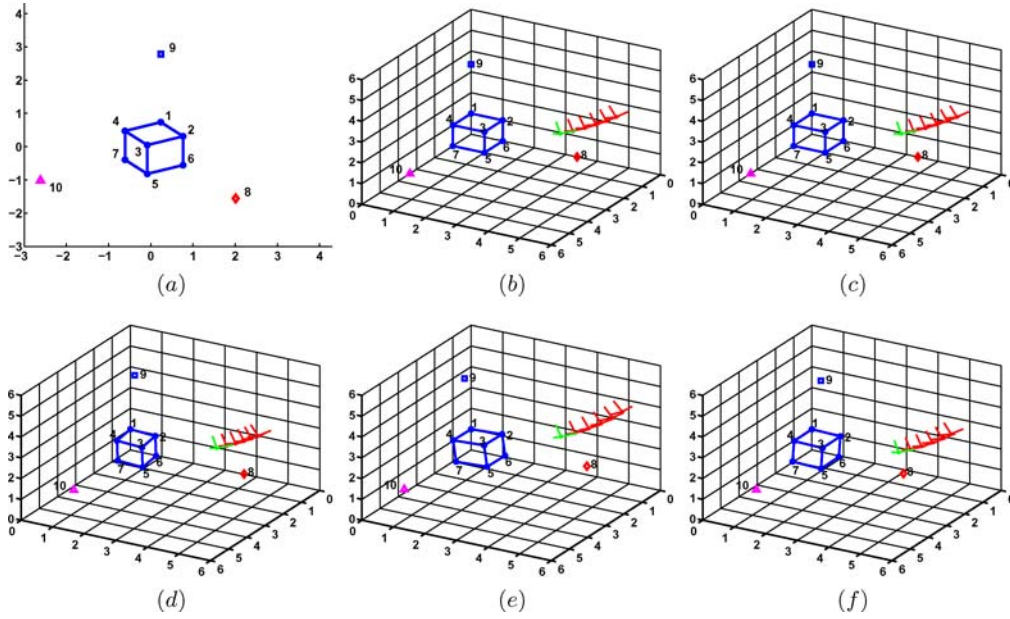


Figure 1. A static cube and 3 points moving along straight lines. (a) Input image. (b) Ground truth 3D shape and camera trajectory. (c) Reconstruction by the closed-form solution. (d) Reconstruction by the method in Bregler et al. (2000). (e) Reconstruction by the method in Brand (2001) after 4000 iterations. (f) Reconstruction by the tri-linear method (Torresani et al., 2001) after 4000 iterations.

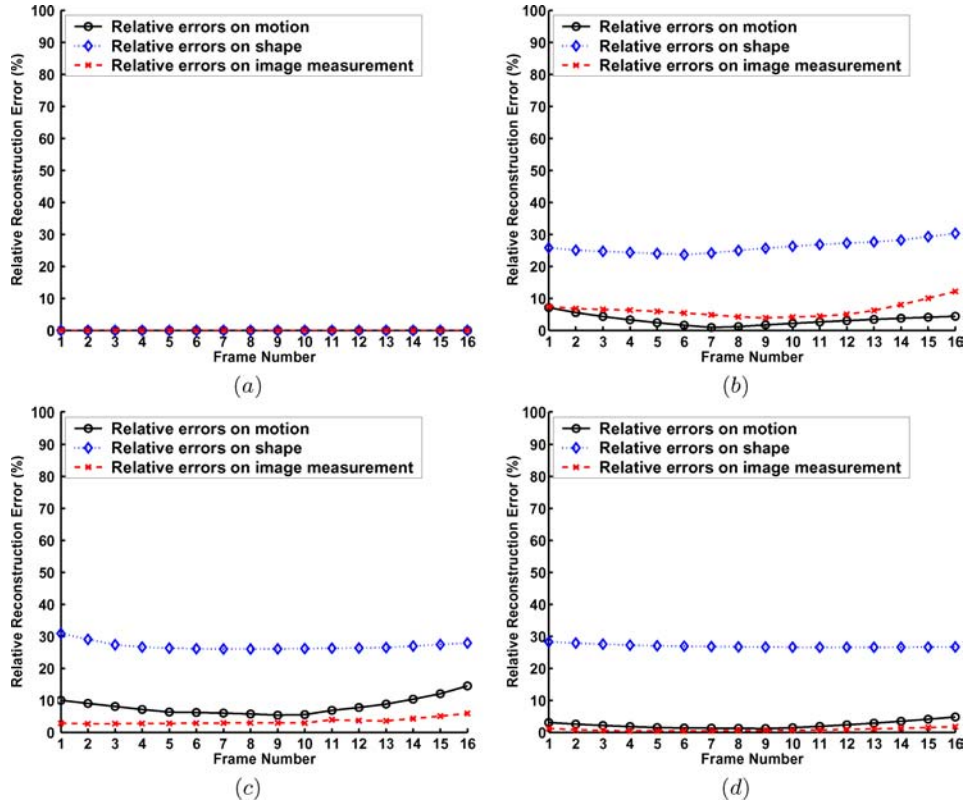


Figure 2. The relative errors on reconstruction of a static cube and 3 points moving along straight lines. (a) By the closed-form solution. (b) By the method in Bregler et al. (2000). (c) By the method in Bregler et al. (2000) after 4000 iterations. (d) By the trilinear method (Torresani et al., 2001) after 4000 iterations. The scaling of the error axis is [0,100%]. Note that our method achieved zero reconstruction errors.

Due to Eq. (6), $M_{2i-1:2i} g_k = c_{ik} R_i$, $i = 1, \dots, F$. Because the rotation matrix R_i is orthonormal, i.e. $\|R_i\| = 1$, we have $R_i = \pm \frac{M_{2i-1:2i} g_k}{\|M_{2i-1:2i} g_k\|}$. The sign of R_i determines which orientations are in front of the camera. It can be either positive or negative, determined by the reference coordinate system. Since g_1, \dots, g_K are estimated independently, they lead to respective rotation sets, each two of which are different up to a 3×3 orthonormal transformation. We choose one set of the rotations to specify the reference coordinate system. Then the signs of the other sets of rotations are determined in such a way that these rotations are consistent with the corresponding references. Finally the orthogonal Procrustes method (Schönemann, 1966) is applied to compute the orthonormal transformations from the rotation sets to the reference. One of the reviewers pointed out that, computing such transformation Φ_k is equivalent to computing a homography at infinity. For details, please refer to Hartley and Zisserman (2000).

The transformed g_1, \dots, g_K form the desired corrective transformation G . The coefficients are then computed by Eq. (6), and the shape bases are recovered by Eq. (5). Their combinations reconstruct the non-rigid 3D shapes.

Note that one could derive the similar basis and rotation constraints directly on $Q = GG^T$, where G is the entire corrective transformation matrix, instead of $Q_k = g_k g_k^T$, where g_k is the k_{th} triple-column of G . It is easy to show that enforcing these constraints also results in a linear closed-form solution of Q . However, the decomposition of Q to recover G is up to a $3K \times 3K$ orthonormal ambiguity transformation. Eliminating this ambiguity alone is as hard as the original problem of determining G .

6. Performance Evaluation

The performance of the closed-form solution was evaluated in a number of experiments.

6.1. Comparison with Three Previous Methods

We first compared the solution with three related methods (Bregler et al., 2000, Brand, 2001, Torresani et al., 2001) in a simple noiseless setting. Figure 1 shows a scene consisting of a static cube and 3 moving points. The measurement included 10 points: 7 visible vertices of the cube and 3 moving points. The 3 points moved along the three axes simultaneously at varying

speed. This setting involved $K = 2$ shape bases, one for the static cube and another for the linear motions. While the points were moving, the camera was rotating around the scene. A sequence of 16 frames were captured. One of them is shown in Fig. 1(a) and (b) demonstrates the ground truth shape in this frame and the ground truth camera trajectory from the first frame till this frame. The three orthogonal green bars show the present camera pose and the red bars display the camera poses in the previous frames. Figure 1(c) to (f) show the structures and camera trajectories reconstructed using the closed-form solution, the method in Bregler et al. (2000), the method in Brand (2001), and the tri-linear method (Torresani et al., 2001), respectively. While the closed-form solution achieved the exact reconstruction with zero error, all the three previous methods resulted in apparent errors, even for such a simple noiseless setting.

Figure 2 demonstrates the reconstruction errors of the four methods on camera rotations, shapes, and image measurements. The error was computed as the percentage relative to the ground truth, $\frac{\|Reconstruction - Truth\|}{\|Truth\|}$. Note that because the space of rotations is a manifold, a better error measurement for rotations is the Riemannian distance, $d(R_i, R_j) = \arccos(\frac{\text{trace}(R_i R_j^T)}{2})$. However it is measured in degrees. For consistency, we used the relative percentage for all the three reconstruction errors.

6.2. Quantitative Evaluation on Synthetic Data

Our approach was then quantitatively evaluated on the synthetic data. We evaluated the accuracy and robustness on three factors: deformation strength, number of shape bases, and noise level. The deformation strength shows how close to rigid the shape is. It is represented by the mean power ratio between each two bases, i.e. $mean_{i,j}(\frac{\max(\|B_i\|, \|B_j\|)}{\min(\|B_i\|, \|B_j\|)})$. Larger ratio means weaker deformation, i.e. the shape is closer to rigid. The number of shape bases represents the flexibility of the shape. A bigger basis number means that the shape is more flexible. Assuming a Gaussian white noise, we represent the noise strength level by the ratio between the Frobenius norm of the noise and the measurement, i.e. $\frac{\|noise\|}{\|W\|}$. In general, when noise exists, a weaker deformation leads to better performance, because some deformation mode is more dominant and the noise relative to the dominant basis is weaker; a bigger basis number results in poorer performance, because the noise relative to each individual basis is stronger.

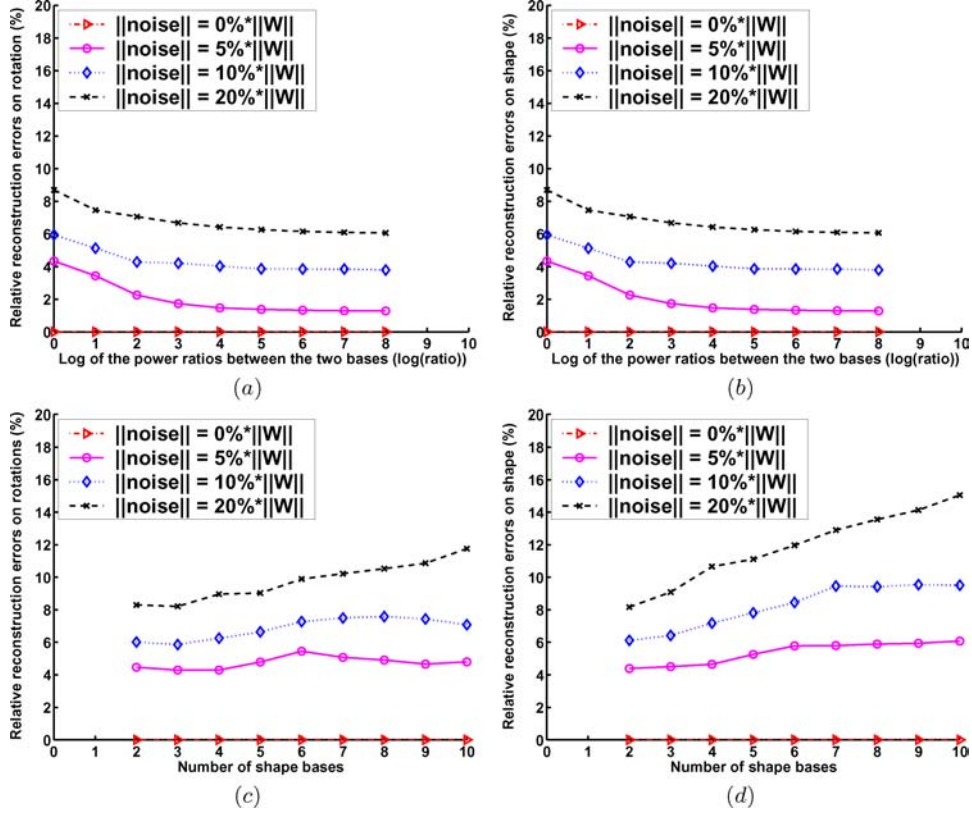


Figure 3. (a) & (b) Reconstruction errors on rotations and shapes under different levels of noise and deformation strength. (c) & (d) Reconstruction errors on rotations and shapes under different levels of noise and various basis numbers. Each curve respectively refers to a noise level. The scaling of the error axis is $[0, 20\%]$.

Figure 3(a) and (b) show the performance of our algorithm under various deformation strength and noise levels on a two bases setting. The power ratios were respectively $2^0, 2^1, \dots$, and 2^8 . Four levels of Gaussian white noise were imposed. Their strength levels were 0, 5, 10, and 20% respectively. We tested a number of trials on each setting and computed the average reconstruction errors on the rotations and 3D shapes. The errors were measured by the relative percentage as in Section 6.1. Figure 3(c) and (d) show the performance of our method under different numbers of shape bases and noise levels. The basis number was 2, 3, \dots , and 10 respectively. The bases had equal powers and thus none of them was dominant. The same noise as in the last experiment was imposed.

In both experiments, when the noise level was 0%, the closed-form solution recovered the exact rotations and shapes with zero error. When there was noise, it achieved reasonable accuracy, e.g. the maximum reconstruction error was less than 15% when the noise

level was 20%. As we expected, under the same noise level, the performance was better when the power ratio was larger and poorer when the basis number was bigger. Note that in all the experiments, the condition number of the linear system consisting of both basis constraints and rotation constraints had order of magnitude $O(10)$ to $O(10^2)$, even if the basis number was big and the deformation was strong. It suggests that our closed-form solution is numerically stable.

6.3. Qualitative Evaluation on Real Video Sequences

Finally we examined our approach qualitatively on a number of real video sequences. One example is shown in Fig. 4. The sequence was taken of an indoor scene by a handheld camera. Three objects, a car, a plane, and a toy person, moved along fixed directions and at varying speeds. The rest of the scene was static. The

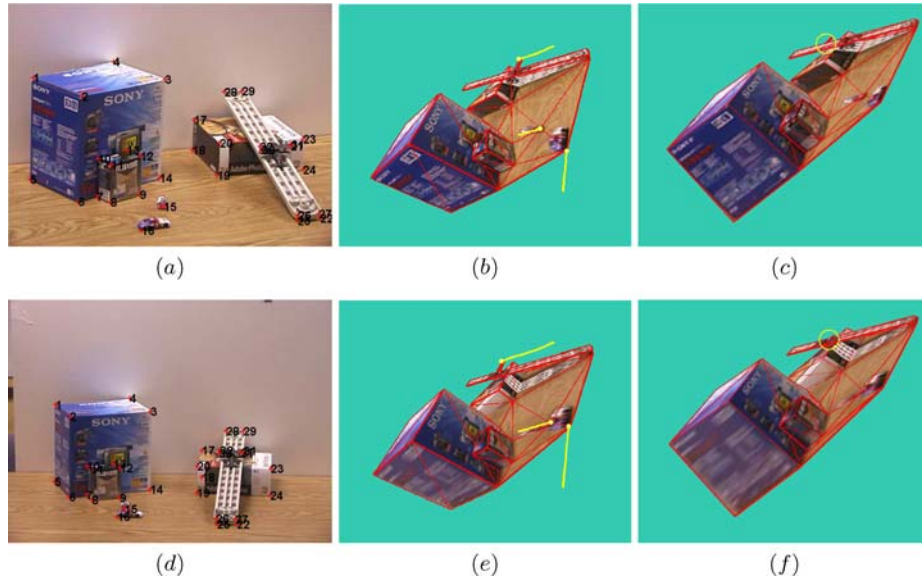


Figure 4. Reconstruction of three moving objects in the static background. (a) & (d) Two input images with marked features. (b) & (e) Reconstruction by the closed-form solution. The yellow lines show the recovered trajectories from the beginning of the sequence until the present frames. (c) & (f) Reconstruction by the method in Brand (2001). The yellow-circled area shows that the plane, which should be on top of the slope, was mistakenly located underneath the slope.

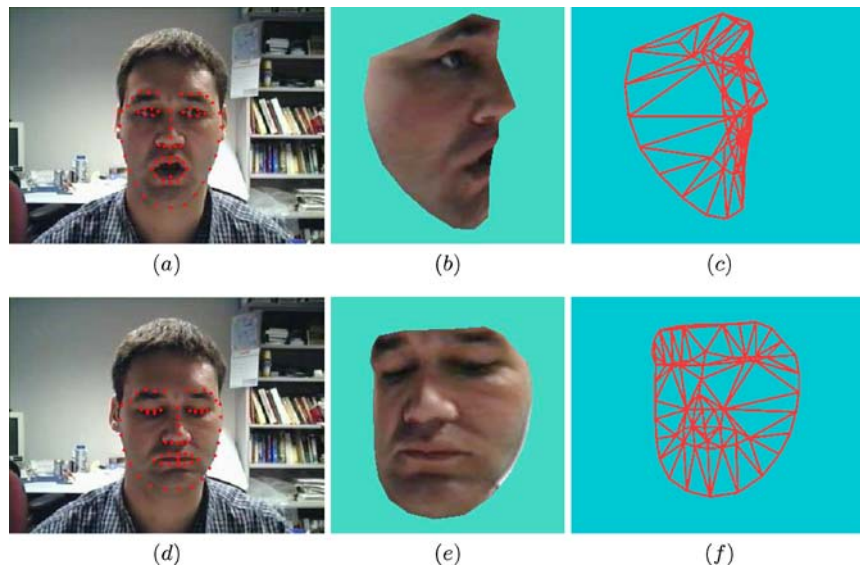


Figure 5. Reconstruction of shapes of human faces carrying expressions. (a) & (d) Input images. (b) & (e) Reconstructed 3D face appearances in novel poses. (c) & (f) Shape wireframes demonstrating the recovered facial deformations such as mouth opening and eye closure.

car and the person moved on the floor and the plane moved along a slope. The scene structure is composed of two bases, one for the static objects and another for the linear motions. 32 feature points tracked across 18

images were used for reconstruction. Two of the them are shown in Fig. 4(a) and (d).

The rank of \tilde{W} is estimated in such a way that 99% of the energy of \tilde{W} remains after the factorization using

the rank constraint. The number of the bases is thus determined by $K = \frac{\text{rank}(\tilde{W})}{3}$. Then the camera rotations and dynamic scene structures were reconstructed using the proposed method. With the recovered shapes, we could view the scene appearance from any novel directions. An example is shown in Fig. 4(b) and (e). The wireframes show the scene shapes and the yellow lines show the trajectories of the moving objects from the beginning of the sequence until the present frames. The reconstruction was consistent with our observation, e.g. the plane moved linearly on top of the slope. Figure 4(c) and (f) show the reconstruction using the method in Brand (2001). The recovered shapes of the boxes were distorted and the plane was incorrectly located underneath the slope, as shown in the yellow circles. Note that occlusion was not taken into account when rendering these images. Thus in the regions that should be occluded, e.g. the area behind the slope, the stretched texture of the occluding objects appeared.

Human faces are highly non-rigid objects and 3D face shapes can be regarded as weighted combinations of certain shape bases that refer to various facial expressions. Thus our approach is capable of reconstructing the deformable 3D face shapes from the 2D image sequence. One example is shown in Fig. 5. The sequence consisted of 236 images that contained facial expressions like eye blinking and mouth opening. 60 feature points were tracked using an efficient 2D Active Appearance Model (AAM) method (Baker and Matthews, 2001). Figure 5(a) and (d) display two of the input images with marked feature points. After reconstructing the shapes and poses, we could view the 3D face appearances in any novel poses. Two examples are shown respectively in Fig. 5(b) and (e). Their corresponding 3D shape wireframes, as shown in Fig. 5(c) and (f), exhibit the recovered facial deformations such as mouth opening and eye closure. Note that the feature correspondences in these experiments were noisy, especially for those features on the sides of the face. The reconstruction performance of our approach demonstrates its robustness to the image noise.

7. Conclusion and Discussion

This paper proposes a linear closed-form solution to the problem of non-rigid shape and motion recovery from a single-camera video. In particular, we have proven that enforcing only the rotation constraints results in ambiguous and invalid solutions. We thus in-

troduce the basis constraints to resolve this ambiguity. We have also proven that imposing both the linear constraints leads to a unique reconstruction of the non-rigid shape and motion. The performance of our algorithm is demonstrated by experiments on both simulated data and real video data. Our algorithm has also been successfully applied to separate the local deformations from the global rotations and translations in the 3D motion capture data (Chai et al., 2003).

Currently our approach does not consider the degenerate deformations. A shape basis is degenerate, if its rank is either 1 or 2, i.e. it limits the shape deformation within a 2D plane. Such planar deformations occur in structures like dynamic scenes or expressive human faces. For example, when a scene consists of several buildings and one car moving straight along a street, the shape basis referring to the rank-1 car translation is degenerate. It is conceivable that, in such degenerate cases, the basis constraints cannot completely resolve the ambiguity of the rotation constraints. We are now exploring how to extend the current method to reconstructing the shapes involving degenerate deformations. Another limitation of our approach is that we assume the weak perspective projection model. It would be interesting to see how the proposed solution could be extended to the full perspective projection model.

Acknowledgments

We would like to thank Simon Baker, Iain Matthews, and Mei Han for providing the image data and feature correspondence used in Section 6.3, thank the ECCV and IJCV reviewers for their valuable comments, and thank Jessica Hodgins for proofreading the paper. Jinxiang Chai was supported by the NSF through EIA0196217 and IIS0205224. Jing Xiao and Takeo Kanade were partly supported by grant R01 MH51435 from the National Institute of Mental Health.

References

- Baker, S. and Matthews, I. 2001. Equivalence and efficiency of image alignment algorithms. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Basclé, B. and Blake, A. 1998. Separability of pose and expression in facial tracing and animation. In *Proc. Int. Conf. Computer Vision*, pp. 323–328.
- Blanz, V. and Vetter, T. 1999. A morphable model for the synthesis of 3D faces. In *Proc. SIGGRAPH'99*, pp. 187–194.

- Brand, M. 2001. Morphable 3D models from video. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Brand, M. and Bhotika, R. 2001. Flexible flow for 3D nonrigid tracking and shape recovery. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Bregler, C., Hertzmann, A. and Biermann, H. 2000. Recovering non-rigid 3D shape from image streams. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Chai, J., Xiao, J., and Hodgins, J. 2003. Vision-based control of 3D facial animation. *Eurographics/ACM Symposium on Computer Animation*.
- Costeira, J. and Kanade, T. 1998. A multibody factorization method for independently moving-objects. *Int. Journal of Computer Vision*, 29(3):159–179.
- Gokturk, S.B., Bouguet, J.Y., and Grzeszczuk, R. 2001. A data driven model for monocular face tracking. In *Proc. Int. Conf. Computer Vision*.
- Han, M. and Kanade, T. 2000. Reconstruction of a scene with multiple linearly moving objects. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Hartley, R. I. and Zisserman, A. 2000. *Multiple View Geometry in Computer Vision*. Cambridge University Press.
- Kanatani, K. 2001. Motion segmentation by subspace separation and model selection. In *Proc. Int. Conf. Computer Vision*.
- Mahamud, S. and Hebert, M. 2000. Iterative projective reconstruction from multiple views. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Poelman, C. and Kanade, T. 1997. A paraperspective factorization method for shape and motion recovery, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 19(3):206–218.
- Schönemann, P.H. 1966. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10.
- Sugaya, Y. and Kanatani, K. 2004. Geometric structure of degeneracy for multi-body motion segmentation. *ECCV Workshop on Statistical Methods in Video Processing*.
- Tomasi, C. and Kanade, T. 1992. Shape and motion from image streams under orthography: A factorization method. *Int. Journal of Computer Vision*, 9(2):137–154.
- Torresani, L., Yang, D., Alexander, G., and Bregler, C. 2001. Tracking and modeling non-rigid objects with rank constraints. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Triggs, B. 1996. Factorization methods for projective structure and motion. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Vidal, R., Soatto, S., Ma, Y., and Sastry, S. 2002. Segmentation of dynamic scenes from the multibody fundamental matrix. *ECCV Workshop on Vision and Modeling of Dynamic Scenes*.
- Vidal, R. and Hartley, R. 2004. Motion segmentation with missing data using power factorization and GPCA. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Wolf, L. and Shashua, A. 2001. Two-body segmentation from two perspective views. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Wolf, L. and Shashua, A. 2002. On projection matrices $P^{\kappa} \rightarrow P^2, \kappa = 3, \dots, 6$, and their applications in computer vision. *Int. J. of Computer Vision*, 48(1):53–67.
- Zelnik-Manor, L. and Irani, M. 2003. Degeneracies, dependencies, and their implications in multi-body and multi-sequence factorizations. In *Proc. Int. Conf. Computer Vision and Pattern Recognition*.
- Zhang, Z., Deriche, R., Faugeras, O., and Luong, Q. 1995. A robust technique for matching two uncalibrated images through the recovery of the unknown epipolar geometry, *Artificial Intelligence*, 78(1/2):87–119.

Copyright of International Journal of Computer Vision is the property of Springer Science & Business Media B.V. and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.