

A Fast Algorithm for Rigid Structure from Image Sequences

Pedro M. Q. Aguiar *

José M. F. Moura

Department of Electrical and Computer Engineering
Carnegie Mellon University, Pittsburgh PA, USA
{aguiar,moura}@ece.cmu.edu

Abstract

The factorization method [1] is a feature-based approach to recover 3D rigid structure from motion. In [2], we extended their framework to recover a parametric description of the 3D shape. In [1, 2], the 3D shape and 3D motion are computed by using an SVD to approximate a matrix that is rank 3 in a noiseless situation. In this paper, we develop a new algorithm that has two relevant advantages over the algorithms of [1, 2]. First, instead of imposing a common origin for the parametric representation of the 3D surface patches, as in [2], we allow the specification of different origins for different patches. This improves the numerical stability of the image motion estimation algorithm and the accuracy of the 3D structure recovery algorithm. Second, we show how to compute the 3D shape and 3D motion by a simple factorization of a modified matrix that is rank 1 in a noiseless situation, instead of a rank 3 matrix as in [1, 2]. This allows the use of very fast algorithms even when using a large number of features (or regions) and large number of frames.

1 Introduction

The factorization method [1] is an elegant method to recover 3D rigid structure from an image sequence. In [1], the 3D positions of the feature points are expressed in terms of cartesian coordinates in a world-centered coordinate system, and the images are modeled as orthographic projections. The 2D projection of each feature point is tracked along the image sequence. The 3D shape and motion are then estimated by factorizing a measurement matrix whose entries are the set of trajectories of the feature point projections. The factorization of the measurement matrix, which is rank 3 in a noiseless situation, is computed by using the Singular Value Decomposition (SVD). The factorization method was extended to

the scaled-orthography and the paraperspective projections in [3].

When the goal is the recovery of a dense representation of the 3D shape, the factorization approach of [1, 3] may not solve the problem satisfactorily because of two drawbacks. First, being feature-based, it would be necessary to track a huge number of features to obtain a dense description of the 3D shape. This is usually impossible because only distinguished points, as brightness corners, can be accurately tracked. Second, even if it is possible to track a large number of features, the computational cost of the SVD involved in the factorization of the measurement matrix would be very high. These drawbacks motivated the extension of the factorization approach to recover a parametric description of the 3D shape, as we did in [2]. Instead of tracking pointwise features, we track regions for which the motion induced on the image plane is described by a single set of parameters.

In this paper, we reformulate the problem. The reformulation leads to a new algorithm that has two relevant advantages over the algorithms of [1, 2]. First, instead of imposing a common origin for the parametric representation of the 3D surface patches, as in [2], we allow the specification of different origins for different patches. This improves the numerical stability of the image motion estimation algorithm and the accuracy of the 3D structure recovery algorithm, as will become clear later. Second, we show how to compute the 3D shape and 3D motion by a simple factorization of a modified matrix that is rank 1 in a noiseless situation, instead of a rank 3 matrix as in [1, 2]. This simplifies the decomposition and normalization stages involved in the factorization approach. We avoid the computation of the SVD by using a fast iterative method to compute the rank 1 matrix that best matches the data. Reference [4] compares the computational cost of the rank 1 factorization with the computational cost of the original rank 3 factorization method [1] for the feature-based case.

*The first author is also affiliated with ISR-IST, Lisboa, Portugal. The first author was partially supported by INVOTAN.

Paper Overview Section 2 describes the scenario. In section 3, we characterize the motion induced in the image plane. Sections 4 and 5 deal with the recovery of the 3D structure from the image motion. In section 6, we detail a real video experiment. Section 7 concludes the paper.

2 Scenario

We consider a rigid body moving in front of the camera. We attach coordinate systems to the object and to the camera. The object coordinate system (o.c.s.) has axes labeled by x, y , and z . The camera coordinate system (c.c.s.) has axes labeled by u, v , and w . We consider that the o.c.s. coincides with the c.c.s. on the first frame. The image plane is defined by the axes u and v . The images are modeled as orthographic projections of the object texture. Our algorithm is easily extended to the scaled-orthography and the paraperspective projections by proceeding as [3] does for the original factorization method.

3D Shape The shape of the rigid object is a parametric description of the object surface. Although our approach is general enough to cope with general parameterizations, we consider in this paper objects whose shape is given by a piecewise planar surface with K patches. The shape parameter vector is $\mathbf{a} = \{a_{00}^k, a_{10}^k, a_{01}^k, 1 \leq k \leq K\}$ where

$$z = a_{00}^k + a_{10}^k(x - x_0^k) + a_{01}^k(y - y_0^k) \quad (1)$$

describes the shape of the patch k in the o.c.s.. In what respects to the representation of the planar patches, the parameters x_0^k and y_0^k can have any value, for example they can be made zero as we did in [2]. In this paper, we allow the specification of general parameters x_0^k, y_0^k . The relevance of this generalization is obvious: the shape of a small patch k with support region $\{(x, y)\}$ located far from the the point (x_0^k, y_0^k) has an high sensibility with respect to the shape parameters a_{10}^k and a_{01}^k . To minimize this sensibility, we choose for (x_0^k, y_0^k) the centroid of the support region of patch k . With this choice, we improve the numerical stability of the image motion estimation algorithm and the accuracy of the 3D structure recovery algorithm.

To simplify the notation, we define the vectors $\mathbf{a}_k = [a_{10}^k, a_{01}^k]^T$, $\mathbf{s} = [x, y]^T$, and $\mathbf{s}_{k0} = [x_0^k, y_0^k]^T$, and rewrite the shape of the patch k as

$$z = a_{00}^k + \mathbf{a}_k^T(\mathbf{s} - \mathbf{s}_{k0}^T). \quad (2)$$

3D Motion We define the 3D motion of the object by specifying the position of the o.c.s. relative to the c.c.s. in terms of $(t_{uf}, t_{vf}, t_{wf}, \Theta_f)$ where (t_{uf}, t_{vf}, t_{wf}) are

the coordinates of the origin of the o.c.s. with respect to the c.c.s. (3D translation), and Θ_f is the rotation matrix that determine the orientation of the o.c.s. relative to the c.c.s. (3D rotation).

A point with coordinates $[x, y, z]^T$ in the o.c.s. has the following coordinates in the c.c.s., at frame f ,

$$\begin{bmatrix} u_f \\ v_f \\ w_f \end{bmatrix} = \begin{bmatrix} i_{xf} & i_{yf} & i_{zf} \\ j_{xf} & j_{yf} & j_{zf} \\ k_{xf} & k_{yf} & k_{zf} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \begin{bmatrix} t_{uf} \\ t_{vf} \\ t_{wf} \end{bmatrix}, \quad (3)$$

where the matrix above is the 3D rotation matrix Θ_f .

3 Image Motion

In this section we show that the motion induced in the image plane by the body-camera 3D motion is affine with different parameterizations for regions corresponding to different patches. We relate the parameters of the affine motion model to the 3D shape and 3D motion parameters.

Under orthography, the point with coordinates (x, y, z) in the o.c.s. projects in frame f to the image point (u_f, v_f) given by

$$\begin{bmatrix} u_f \\ v_f \end{bmatrix} = \mathbf{M}_f \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \mathbf{t}_f, \quad (4)$$

where \mathbf{M}_f collects the first and second rows of the 3D rotation matrix introduced in expression (3) and $\mathbf{t}_f = [t_{uf}, t_{vf}]^T$.

Consider a generic point in the object surface with coordinates $\mathbf{s} = [x, y]^T$ and z given by expression (2). We denote by $\mathbf{u}_f(\mathbf{s}) = [u_f(\mathbf{s}), v_f(\mathbf{s})]^T$ the trajectory of the projection of the point \mathbf{s} in the image plane. Since we have chosen the coordinate systems to coincide on the first frame, we have $\mathbf{u}_1(\mathbf{s}) = \mathbf{s}$. At frame f , the point \mathbf{s} projects according to expression (4), to the image point

$$\mathbf{u}_f(\mathbf{s}) = \mathbf{N}_f \mathbf{s} + \mathbf{n}_f z + \mathbf{t}_f, \quad (5)$$

where we have decomposed the matrix \mathbf{M}_f as $\mathbf{M}_f = [\mathbf{N}_f, \mathbf{n}_f]$ where \mathbf{N}_f collects the first and second columns of \mathbf{M}_f and \mathbf{n}_f is the third column of \mathbf{M}_f .

Affine Motion Model By inserting expression (2) into expression (5), we express the image displacement between frame 1 and frame f in terms of the 3D shape and 3D motion parameters, for the points \mathbf{s} that fall into patch k of the object surface. After simple manipulations, we obtain

$$\mathbf{u}_f^k(\mathbf{s}) = (\mathbf{N}_f + \mathbf{n}_f \mathbf{a}_k^T)(\mathbf{s} - \mathbf{s}_{k0}^T) + \mathbf{N}_f \mathbf{s}_{k0}^T + \mathbf{n}_f a_{00}^k + \mathbf{t}_f. \quad (6)$$

Denoting the matrix that multiplies $(s - s_0^k)$ and the vector corresponding to the term independent of s by

$$\begin{cases} D_f^k = N_f + n_f a_k^T \\ d_f^k = N_f s_0^k + n_f a_{00}^k + t_f \end{cases}, \quad (7)$$

we rewrite expression (6) as

$$u_f^k(s) = D_f^k (s - s_0^k) + d_f^k. \quad (8)$$

Expression (8) shows that the image coordinates at frame f , u_f , of the the points belonging to the object surface are affine mappings of their image coordinates frame 1, $u_1 = s$. Expression (7) relates the coefficients of the affine motion models for each patch k to the 3D motion parameters and the 3D shape parameters corresponding to patch k .

Image Motion Estimation Except for particular 3D motions, the image motion corresponding to different surface patches is described by different affine parameterizations. The problem of estimating the support regions of the surface patches leads to the segmentation of the image motion field. The segmentation according to image motion has been widely addressed in the past, see for example [5, 6]. We use the simple method of sliding a rectangular window across the image and detect abrupt changes in the affine motion parameters.

Another possible way to use our structure from motion approach is to select *a priori* the support regions of the surface patches. In fact, our framework is general enough to accommodate the feature tracking approach because it corresponds to selecting *a priori* a set of small (pointwise) support regions with shape described by $z = \text{constant}$ in each region. In reference [4] we exploit the feature-based approach.

4 Rigid Structure from Motion

The problem of inferring 3D rigid structure from the image motion is formulated as estimating the 3D motion parameters $\{N_f, n_f, t_f, 2 \leq f \leq F\}$ and the 3D shape parameters $\{a_{00}^k, a_k, 1 \leq k \leq K\}$ from the image motion parameters $\{D_f^k, d_f^k, 2 \leq f \leq F, 1 \leq k \leq K\}$ by inverting the overconstrained set of equations of expression (7).

We start by estimating the translation. By choosing the object coordinate system in such a way that $\sum_k a_{00}^k = 0$ and the image origin in such a way that $\sum_k s_0^k = [0, 0]^T$, we obtain the Least Squares (LS) estimate for the translation vector t_f as the mean of the vectors $\{d_f^k, 1 \leq k \leq K\}$,

$$\hat{t}_f = \frac{1}{K} \sum_{k=1}^K d_f^k. \quad (9)$$

To eliminate the dependence of the image motion parameters on the translation, we replace the translation estimates into expression (7) and define a new set of parameters $\{\tilde{d}_f^k\}$ related to $\{d_f^k\}$ by

$$\tilde{d}_f^k = d_f^k - \frac{1}{K} \sum_{l=1}^K d_f^l. \quad (10)$$

Defining the matrices R_f^k and S_k^T as

$$R_f^k = \begin{bmatrix} D_f^k & \tilde{d}_f^k \end{bmatrix} \text{ and } S_k^T = \begin{bmatrix} I_{2 \times 2} & s_0^k \\ a_k^T & a_{00}^k \end{bmatrix}, \quad (11)$$

we rewrite the equation system (7) in matrix format as

$$R_f^k = M_f S_k^T. \quad (12)$$

Expression (12) relates the image motion parameters at frame f and patch k to the 3D rotation at frame f and the 3D shape parameters for the patch k .

To make explicit the entire set of equations that arise from considering every patch $1 \leq k \leq K$ and every frame $2 \leq f \leq F$, we define the $2(F-1) \times 3K$ matrix R of image motion parameters, the $2(F-1) \times 3$ matrix M of 3D rotation parameters, and the $3K \times 3$ matrix S of 3D shape parameters as

$$R = \begin{bmatrix} R_2^1 & \cdots & R_2^K \\ \vdots & \ddots & \vdots \\ R_F^1 & \cdots & R_F^K \end{bmatrix}, M = \begin{bmatrix} M_2 \\ \vdots \\ M_F \end{bmatrix}, S = \begin{bmatrix} S_1 \\ \vdots \\ S_K \end{bmatrix}, \quad (13)$$

and we write the relation between the image motion parameters and the 3D structure parameters as

$$R = M S^T. \quad (14)$$

The matrix R of image motion parameters is highly rank deficient. In a noiseless situation, R is rank 3 reflecting the high redundancy in the data, due to the rigidity of the object.

5 Rank 1 Factorization

Estimating the 3D shape and 3D rotation parameters given the observation matrix R is a nonlinear LS problem. The factorization approach [1, 2] finds a sub-optimal solution to this problem in two stages. The first stage, *decomposition stage*, solves $R = M S^T$ in the LS sense by computing the SVD of the matrix R and selecting the 3 largest singular values. From $R \simeq U \Sigma V^T$, the solution is $M = U \Sigma^{\frac{1}{2}} A$ and $S^T = A^{-1} \Sigma^{\frac{1}{2}} V^T$ where A is a non-singular 3×3 matrix. The second stage, *normalization stage*, computes A by approximating the constraints imposed by the structure of the matrices M and S .

The formulation we adopt in this paper takes advantage of the fact that the first two rows of \mathbf{S}^T are known. The problem is then reduced to, given \mathbf{R} , compute \mathbf{M} and $\{a_{mn}^k\}$. We also perform in sequence the decomposition and normalization stages. These stages are as follows. For more details see [4].

Decomposition Because the first two rows of \mathbf{S}^T are known (see expressions (11), and (13)), we show that the unconstrained bilinear problem $\mathbf{R} = \mathbf{M}\mathbf{S}^T$ is solved by the factorization of a rank 1 matrix, rather than a rank 3 matrix like in [1, 2]. Define $\mathbf{M} = [\mathbf{M}_0, \mathbf{m}_3]$ and $\mathbf{S} = [\mathbf{S}_0, \mathbf{a}]$. Matrices \mathbf{M}_0 and \mathbf{S}_0 contain the first two columns of \mathbf{M} and \mathbf{S} , respectively, \mathbf{m}_3 is the third column of \mathbf{M} , and \mathbf{a} is the third column of \mathbf{S} . We decompose the shape parameter vector \mathbf{a} into the component that belongs to the space spanned by the columns of \mathbf{S}_0 and the component orthogonal to this space as $\mathbf{a} = \mathbf{S}_0\mathbf{b} + \mathbf{a}_1$, with $\mathbf{a}_1^T \mathbf{S}_0 = [0, 0]$. Using these definitions, we rewrite \mathbf{R} as

$$\mathbf{R} = \mathbf{M}_0\mathbf{S}_0^T + \mathbf{m}_3\mathbf{b}^T\mathbf{S}_0^T + \mathbf{m}_3\mathbf{a}_1^T. \quad (15)$$

The decomposition stage is formulated as

$$\min_{\mathbf{M}_0, \mathbf{m}_3, \mathbf{b}, \mathbf{a}_1} \left\| \mathbf{R} - \mathbf{M}_0\mathbf{S}_0^T - \mathbf{m}_3\mathbf{b}^T\mathbf{S}_0^T - \mathbf{m}_3\mathbf{a}_1^T \right\|_F, \quad (16)$$

where $\|\cdot\|_F$ denotes the Frobenius norm. By solving the linear LS for \mathbf{M}_0 in terms of the other variables, we get

$$\widehat{\mathbf{M}}_0 = \mathbf{R}\mathbf{S}_0 \left(\mathbf{S}_0^T \mathbf{S}_0 \right)^{-1} - \mathbf{m}_3\mathbf{b}^T, \quad (17)$$

where we used the orthogonality between \mathbf{a}_1 and \mathbf{S}_0 . By replacing $\widehat{\mathbf{M}}_0$ in (16), we get

$$\min_{\mathbf{m}_3, \mathbf{a}_1} \left\| \widetilde{\mathbf{R}} - \mathbf{m}_3\mathbf{a}_1^T \right\|_F, \quad (18)$$

$$\text{where } \widetilde{\mathbf{R}} = \mathbf{R} \left[\mathbf{I} - \mathbf{S}_0 \left(\mathbf{S}_0^T \mathbf{S}_0 \right)^{-1} \mathbf{S}_0^T \right]. \quad (19)$$

We see that the decomposition stage does not determine the vector \mathbf{b} . This is because the component of \mathbf{a} that lives in the space spanned by the columns of \mathbf{S}_0 does not affect the space spanned by the columns of the entire matrix \mathbf{S} and the decomposition stage restricts only this last space.

The solution for \mathbf{m}_3 and \mathbf{a}_1 is given by the rank 1 matrix that best approximates $\widetilde{\mathbf{R}}$. In a noiseless situation, $\widetilde{\mathbf{R}}$ is rank 1 (we get $\widetilde{\mathbf{R}} = \mathbf{m}_3\mathbf{a}_1^T$ by replacing (15) in (19)). By computing the largest singular value of $\widetilde{\mathbf{R}}$ and the associated singular vectors, we get

$$\widetilde{\mathbf{R}} \simeq \mathbf{u}\sigma\mathbf{v}^T, \quad \widehat{\mathbf{m}}_3 = \alpha\mathbf{u}, \quad \widehat{\mathbf{a}}_1^T = \frac{\sigma}{\alpha}\mathbf{v}^T, \quad (20)$$

where α is a normalizing scalar different from 0. To compute \mathbf{u} , σ , and \mathbf{v} we use a fast algorithm outlined in [4]. This makes our decomposition stage simpler than the one in [1, 2]. In fact, $\widetilde{\mathbf{R}}$ in (19) is \mathbf{R} multiplied by the orthogonal projector onto the orthogonal complement of the space spanned by the columns of \mathbf{S}_0 . This projection reduces the rank of the problem from 3 (matrix \mathbf{R}) to 1 (matrix $\widetilde{\mathbf{R}}$).

Normalization We compute α and \mathbf{b} by imposing the constraints that come from the structure of \mathbf{M} . By replacing $\widehat{\mathbf{m}}_3$ in (17), we get

$$\widehat{\mathbf{M}} = \begin{bmatrix} \widehat{\mathbf{M}}_0 & \widehat{\mathbf{m}}_3 \end{bmatrix} = \mathbf{N} \begin{bmatrix} \mathbf{I}_{2 \times 2} & \mathbf{0}_{2 \times 1} \\ -\alpha\mathbf{b}^T & \alpha \end{bmatrix}, \quad (21)$$

$$\text{where } \mathbf{N} = \begin{bmatrix} \mathbf{R}\mathbf{S}_0 \left(\mathbf{S}_0^T \mathbf{S}_0 \right)^{-1} & \mathbf{u} \end{bmatrix}. \quad (22)$$

The constraints imposed by the structure of \mathbf{M} are the unit norm of each row and the orthogonality between row $2j$ and row $2j - 1$. In terms of \mathbf{N} , α , and \mathbf{b} , the constraints are

$$\mathbf{n}_i^T \begin{bmatrix} \mathbf{I}_{2 \times 2} & -\alpha\mathbf{b} \\ -\alpha\mathbf{b}^T & \alpha^2(1 + \mathbf{b}^T\mathbf{b}) \end{bmatrix} \mathbf{n}_i = 1 \quad \text{and} \quad (23)$$

$$\mathbf{n}_{2j}^T \begin{bmatrix} \mathbf{I}_{2 \times 2} & -\alpha\mathbf{b} \\ -\alpha\mathbf{b}^T & \alpha^2(1 + \mathbf{b}^T\mathbf{b}) \end{bmatrix} \mathbf{n}_{2j-1} = 0, \quad (24)$$

where \mathbf{n}_i^T denotes the row i of matrix \mathbf{N} . We compute α and \mathbf{b} from the linear LS solution of the system above in a similar way to the one described in [1]. The normalization stage is also simpler than the one in [1] because the number of unknowns is 3 (α and $\mathbf{b} = [b_1, b_2]^T$) as opposed to the 9 entries of a generic 3×3 normalization matrix.

6 Experiment

We used a hand hold taped video sequence of 30 frames showing a box over a carpet. Figure 1 shows frames 1 and 10 of the video sequence. The 3D shape of the scene is well described in terms of four planar patches. One corresponds to the floor, and the other three correspond to the three visible faces of the box. The camera motion was approximately a rotation around the box.

We processed the box video sequence by using the method described above. We start by estimating the affine motion parameters. The plots in figure 2 represent the time evolution of the affine motion parameters. The 6 affine motion parameters are the entries of the 2×2 matrix \mathbf{D}_f^k and the 2×1 vector \mathbf{d}_f^k , see expression (8). The top four plots of figure 2 represent the entries of \mathbf{D}_f^k as a function of f for each of the four

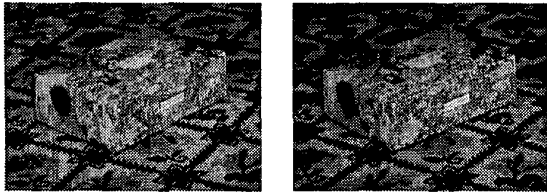


Figure 1: Frames 1 and 10 of the box video sequence.

planar patches. The bottom two plots represent d_f^k . We used four different line types to identify each of the planar patches. The solid line corresponds to patch 1 (the left side vertical face of the box in the frames of figure 1). The dotted line corresponds to patch 2 (the right side vertical face of the box). The dash-dotted line corresponds to patch 3 (the top of the box). The dashed line corresponds to patch 4 (the floor). We see the evolution of the set of affine parameters is distinct for each surface patch, in particular see the evolution of D_{11} , D_{12} , and d_1 .

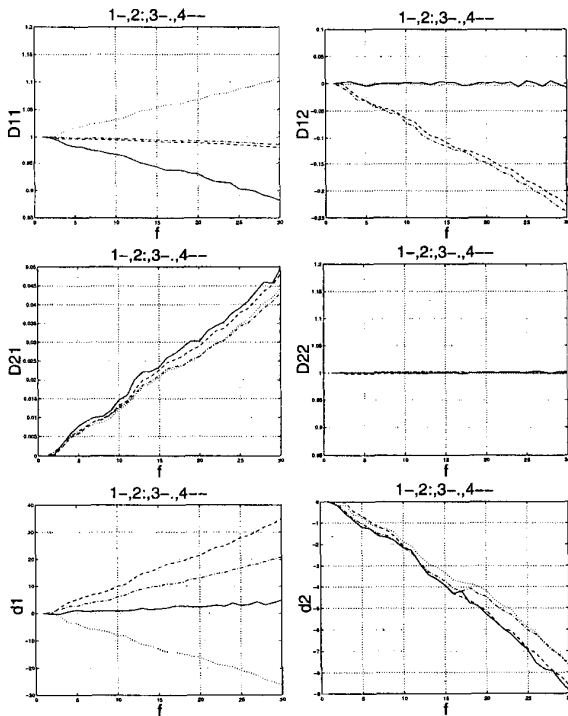


Figure 2: Estimates of the image motion parameters.

From the affine motion parameters of figure 2, we have recovered the 3D structure of the scene by using

the rank 1 factorization method. Figure 3 shows two perspective views of the reconstructed 3D shape with the scene texture mapped on it. We see that the angles between the planar patches are correctly recovered.

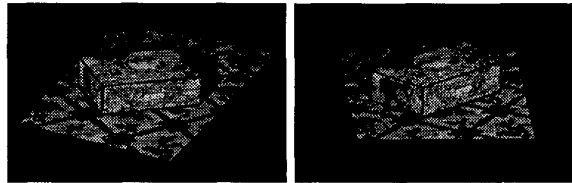


Figure 3: Two perspective views of the reconstructed 3D shape and texture.

7 Conclusion

We proposed a fast method to recover 3D rigid structure from motion. Rather than relying on the tracking of pointwise features, the image motion estimation step makes use of the affine motion parameterization for larger regions, leading to robust estimates of the image motion parameters. The 3D structure from motion step is robust because it takes into account the rigidity of scene over a set of frames. This step is accomplished with very low computational cost by factorizing a rank 1 matrix. The experimental results illustrate the performance of the method.

References

- [1] Carlo Tomasi and Takeo Kanade. Shape and motion from image streams under orthography: a factorization method. *IJCV*, 9(2):137–154, 1992.
- [2] Pedro M. Q. Aguiar and José M. F. Moura. Video representation via 3D shaped mosaics. In *IEEE ICIP*, Chicago, USA, October 1998.
- [3] Conrad J. Poelman and Takeo Kanade. A paraperspective factorization method for shape and motion recovery. *IEEE PAMI*, 19(3):206–218, 1997.
- [4] Pedro M. Q. Aguiar and José M. F. Moura. Factorization as a rank 1 problem. In *IEEE CVPR*, Fort Collins, USA, June 1999.
- [5] H. Zheng and S. Blostein. Motion-based object segmentation and estimation using the MDL principle. *IEEE Trans. Image Processing*, 4(9), 1995.
- [6] M. M. Chang, A. M. Tekalp, and M. I. Sezan. Simultaneous motion estimation and segmentation. *IEEE Trans. Image Processing*, 6(9), 1997.