

# Computing Smooth Time Trajectories for Camera and Deformable Shape in Structure from Motion with Occlusion

Paulo F.U. Gotardo, *Member, IEEE*, and Aleix M. Martinez, *Member, IEEE*

**Abstract**—We address the classical computer vision problems of rigid and nonrigid structure from motion (SFM) with occlusion. We assume that the columns of the input observation matrix  $\mathbf{W}$  describe smooth 2D point trajectories over time. We then derive a family of efficient methods that estimate the column space of  $\mathbf{W}$  using compact parameterizations in the Discrete Cosine Transform (DCT) domain. Our methods tolerate high percentages of missing data and incorporate new models for the smooth time trajectories of 2D-points, affine and weak-perspective cameras, and 3D deformable shape. We solve a rigid SFM problem by estimating the smooth time trajectory of a single camera moving around the structure of interest. By considering a weak-perspective camera model from the outset, we directly compute euclidean 3D shape reconstructions without requiring postprocessing steps such as euclidean upgrade and bundle adjustment. Our results on real SFM data sets with high percentages of missing data compared positively to those in the literature. In nonrigid SFM, we propose a novel *3D shape trajectory* approach that solves for the deformable structure as the smooth time trajectory of a single point in a linear shape space. A key result shows that, compared to state-of-the-art algorithms, our nonrigid SFM method can better model complex articulated deformation with higher frequency DCT components while still maintaining the low-rank factorization constraint. Finally, we also offer an approach for nonrigid SFM when  $\mathbf{W}$  is presented with missing data.

**Index Terms**—Structure from motion, matrix factorization, missing data, camera trajectory, shape trajectory.

## 1 INTRODUCTION

ACCURATELY describing a data matrix as a product of two low-rank factors, matrix factorization, is a fundamental task in computer vision and pattern recognition. This paper focuses on the classical computer vision problem of matrix factorization in rigid and nonrigid structure from motion (SFM) [5], [29]. The goal in SFM is to jointly estimate the 3D scene structure and relative camera motion from corresponding 2D points in a sequence of images. Applications of SFM include autonomous navigation, image augmentation, and the construction of rigid and deformable 3D models from images [11], [15]. The modeling of deformable shapes such as the human hand, face, and body is also of particular importance in computer graphics, and human-computer interaction (e.g., [10]).

While techniques for rigid SFM have matured considerably over the past two decades [6], [7], [8], [12], [16], [17], [19], [21], [28], nonrigid SFM is still a very difficult problem, especially for complex articulated deformations [30]. The difficulty in providing good solutions reflects the under-constrained nature of SFM once the rigidity assumption is removed. Recent research has thus focused on the definition of new constraints (priors) to solve this problem [1], [2], [3], [23], [24], [26], [30], [32], [33], [34].

In the standard matrix factorization approach to SFM [5], [29], each column of the input matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$  has a sequence of 2D coordinates of the same 3D structure point as observed from different locations. Considering  $\mathbf{W}$  of a predefined low-rank  $r \leq \min(m, n)$ , the SFM solution is obtained from the factorization

$$\mathbf{W} = \mathbf{M}\mathbf{S}, \quad \mathbf{M} \in \mathbb{R}^{m \times r}, \quad \mathbf{S} \in \mathbb{R}^{r \times n}. \quad (1)$$

In rigid SFM,  $r = 4$  and  $\mathbf{S}$  describes the 3D shape observed by the cameras in  $\mathbf{M}$ . In nonrigid SFM, each observed shape of the deformable structure is represented in a linear shape space defined by  $K$  basis shapes in  $\mathbf{S}$ . Then,  $r = 3K + 1$  and  $\mathbf{M}$  includes the cameras and also the shape coordinates in terms of basis  $\mathbf{S}$ . In both cases, factors  $\mathbf{M}$  and  $\mathbf{S}$  may be obtained from the singular value decomposition (SVD) of  $\mathbf{W}$ . In practice, however, a large portion of the 2D observations in  $\mathbf{W}$  is often missing because of occlusions. Therefore, standard matrix factorization algorithms such as SVD [13] cannot be directly used. To solve SFM, it is also necessary to overcome other challenges such as tracking errors [19] and degeneracies in the assumed camera motion and shape deformation—i.e., when the above constraints on the rank  $r$  of  $\mathbf{W}$  do not hold [21], [32], [34].

To compute more accurate and efficient solutions to rigid and nonrigid SFM with occlusion, we start by assuming that each column of  $\mathbf{W}$  represents the smooth time trajectory of a 2D point. Such 2D trajectories are usually provided by a feature tracking algorithm that operates on a monocular video sequence. Equivalently, we assume that the 2D observations were obtained by a single camera moving smoothly around the structure of interest (Fig. 1). Our factorization approach tolerates missing data and defines

- The authors are with the Department of Electrical and Computer Engineering, The Ohio State University, 2015 Neil Avenue, 205 Dreese Laboratories, Columbus, OH 43210. E-mail: {gotardop, aleix}@ece.osu.edu.

Manuscript received 22 Feb. 2010; revised 19 Oct. 2010; accepted 28 Jan. 2011; published online 1 Mar. 2011.

Recommended for acceptance by S. Belongie.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2010-02-0119.

Digital Object Identifier no. 10.1109/TPAMI.2011.50.

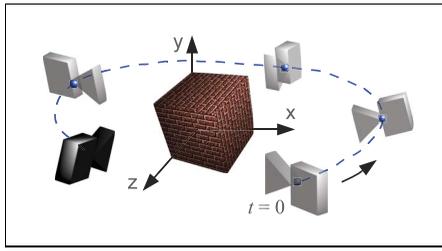


Fig. 1. The smooth trajectory (dashed line) of a single camera moving around the structure of interest (cube) over time ( $t$ ).

constraints that considerably reduce the number of unknowns that need to be estimated. First, we consider factor  $\mathbf{S}$  only implicitly, reformulating the factorization problem in terms of  $\mathbf{M}$  alone. Second, because  $\mathbf{M}$  is a basis for the smooth 2D trajectories in the columns of  $\mathbf{W}$ , we solve for a compact representation of  $\mathbf{M}$  in a subspace of the Discrete Cosine Transform (DCT) basis vectors. Our algorithms are remarkably efficient in cases of long image sequences and many imaged points (i.e., when  $\mathbf{W}$  has high-dimensional column and row spaces).

We solve for rigid SFM by estimating only the smooth time trajectory of a camera's projection plane, Fig. 1. Also, by considering the weak-perspective camera model from the outset, we directly solve for euclidean shape and cameras. Therefore, our method does not require post-processing algorithms that upgrade and refine the initial affine solution (e.g., bundle adjustment [15]).

In nonrigid SFM, we assume that the smooth 2D trajectories in  $\mathbf{W}$  also reflect the smooth deformation of the observed 3D structure over time. The deformable structure is represented as a single point moving smoothly in a linear shape space, Fig. 2. We then solve for a single, smooth 3D shape trajectory with time coordinates in factor  $\mathbf{M}$ . Thus, the associated basis shapes are defined only implicitly in factor  $\mathbf{S}$ . Once  $\mathbf{M}$  has been estimated,  $\mathbf{S}$  is then trivially computed from  $\mathbf{M}$  and  $\mathbf{W}$  by solving simple linear systems of equations.

Our work is most closely related to the recent nonrigid SFM method by Akhter et al. [2]. They propose a factorization approach that does not define a linear shape space, but instead recovers *independent 3D point trajectories* over time. The DCT vectors are used as a basis for individual 3D point trajectories. Their method has provided some of the best results on highly articulated shapes to date. However, the main problem is the method's inability to use additional, higher frequency DCT vectors without increasing the rank of the resulting matrix factors. Thus, its application is restricted to structures with slow and smooth deformation. The method does not address the missing data problem either;  $\mathbf{W}$  is assumed complete and factorized using SVD.

In this paper, we also present an interpretation of Akhter et al.'s approach in terms of a linear shape space and show that their results correspond to coarse approximations of our solutions. Our nonrigid SFM method provides better results on complex articulated deformations due to its more effective use of higher frequency DCT components without increasing the factorization rank. Furthermore, we use the DCT basis to model a camera's trajectory and efficiently solve for rigid SFM. Finally, we also contribute with missing data approaches for both rigid and nonrigid SFM.

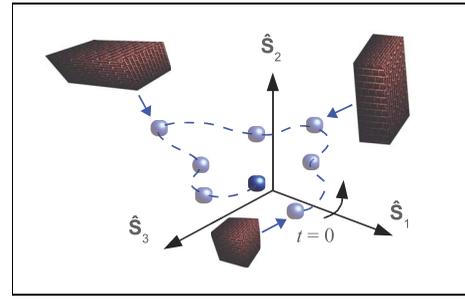


Fig. 2. The smooth time trajectory (dashed line) of a deforming 3D shape (cube). The object shape is represented by a single point (blue sphere) with coordinates relative to 3D basis shapes  $\hat{\mathbf{S}}_1$ ,  $\hat{\mathbf{S}}_2$ , and  $\hat{\mathbf{S}}_3$  (not shown). As the shape deforms smoothly over time ( $t$ ), its point representation describes a single, smooth 3D shape trajectory.

One main group of related methods that address the missing data problem are known as *batch algorithms* [14], [17], [19], [23], [28]. These methods propose strategies for combining partial rank- $r$  factorizations obtained for complete subblocks of  $\mathbf{W}$ . Examples include the subspace constraint algorithms that reconstruct  $\mathbf{W}$  by first building its row null-space [17], [19], column null-space [14], [23], or one of its range spaces [28]. The main problem with these methods is their suboptimality, since errors in the factorization of subblocks are propagated to the subsequent optimization stage. Also, the rank- $r$  constraint may not apply to all subblocks in case of degenerate motion and deformation [21], [23].

A second group of missing data approaches include iterative methods that use all data at once without searching for complete subblocks in  $\mathbf{W}$  [6], [8], [16], [24], [30], [31]. Among these, alternation methods [16], [24], [30] iteratively solve for subsets of unknowns while the others remain fixed. For instance, PowerFactorization (PF) [16] solves for factors  $\mathbf{M}$  and  $\mathbf{S}$  in a simple, alternated least squares manner. It presents very slow convergence when a considerable amount of data is missing. However, the method is useful in initializing faster Newton methods [6] that minimize matrix fitting error in terms of  $\mathbf{M}$  and  $\mathbf{S}$ . The Levenberg-Marquardt-Subspace (LM-S) method in [8] considers  $\mathbf{S}$  as an implicit function of  $\mathbf{M}$  and  $\mathbf{W}$  and solves for  $\mathbf{M}$  only. Despite its superior performance, proper initialization of this method remains an open problem. The method's complexity also makes it difficult to integrate additional constraints into the factorization problem.

In the following, we present our *Column Space Fitting* (CSF) method that computes a rank- $r$  basis  $\mathbf{M}$  for the column space of a matrix with missing data. Using a simple Gauss-Newton-based approximation to the Hessian matrix, our method provides equivalent or better solutions for  $\mathbf{M}$  when compared to LM-S on matrices with high percentage of missing data. Furthermore, the simplicity of our approach allows us to easily consider additional constraints in the factorization problem, such as a mean column vector and a predefined reference basis for  $\mathbf{M}$ . We then solve for a compact representation of  $\mathbf{M}$  in a subspace defined in terms of the DCT basis vectors. As a result, we offer a family of CSF algorithms for rigid and nonrigid SFM problems. Our methods always start from a deterministic initialization corresponding to a "coarse" solution for  $\mathbf{M}$ .

This paper is organized as follows: Section 2 derives our general CSF approach for matrix factorization with missing

data. In Section 3, we solve for rigid SFM by estimating a smooth time trajectory of an affine or a weak-perspective camera. Section 4 describes our 3D shape trajectory approach to nonrigid SFM. Sections 5 and 6 present experimental results and conclusion.

## 2 SOLVING FOR THE COLUMN SPACE OF A MATRIX WITH MISSING DATA

We formulate the rank- $r$  factorization of incomplete matrix  $\mathbf{W}$  as solving exclusively for a column space basis  $\mathbf{M}$  of rank- $r$ . To facilitate the introduction of additional constraints into the factorization procedure, especially in SFM problems, we first derive our Column Space Fitting algorithm. We then show how to solve for  $\mathbf{M}$  in the subspace spanned by a predefined bases. When such a basis is defined in terms of the DCT basis vectors, we propose a simple, deterministic initial form of  $\mathbf{M}$ .

### 2.1 General CSF Approach

Consider  $\mathbf{W} \in \mathbb{R}^{m \times n}$  as in (1) and note the related equation  $\mathbf{S} = \mathbf{M}^\dagger \mathbf{W}$ , where  $\dagger$  denotes the Moore-Penrose pseudo-inverse [13].

Assuming  $\mathbf{W}$  is presented with missing data, let the complete vector  $\mathbf{w}_j \in \mathbb{R}^{m_j}$  ( $m_j \leq m$ ) denote all of the observed entries in the  $j$ th column of  $\mathbf{W}$ . Also, define  $\mathbf{\Pi}_j \in \mathbb{R}^{m_j \times m}$  as a row-amputated identity matrix such that  $\mathbf{M}_j = \mathbf{\Pi}_j \mathbf{M}$  has the rows in  $\mathbf{M}$  that correspond to the rows of entries in  $\mathbf{w}_j$ . Then,  $\mathbf{w}_j = \mathbf{M}_j \mathbf{s}_j$ , with the complete  $j$ th column of  $\mathbf{S}$  defined as  $\mathbf{s}_j = \mathbf{M}_j^\dagger \mathbf{w}_j \in \mathbb{R}^r$ .

The goal is to minimize

$$f(\mathbf{M}) = \frac{1}{2} \sum_j \left\| (\mathbf{I} - \mathbf{M}_j \mathbf{M}_j^\dagger) \mathbf{w}_j \right\|_F^2, \quad (2)$$

where  $\|\cdot\|_F$  is the Frobenius norm. To better understand the meaning of (2), let the projector on the orthogonal space of  $\mathbf{M}_j$  be

$$\mathbf{P}_j^\perp = (\mathbf{I} - \mathbf{M}_j \mathbf{M}_j^\dagger) \in \mathbb{R}^{m_j \times m_j}. \quad (3)$$

Rewriting (2) in terms of residual (error) vectors  $\mathbf{r}_j \in \mathbb{R}^{m_j}$ , we have

$$f(\mathbf{M}) = \frac{1}{2} \sum_j \mathbf{r}_j^T \mathbf{r}_j, \quad \mathbf{r}_j = \mathbf{P}_j^\perp \mathbf{w}_j. \quad (4)$$

Therefore, we minimize the overall matrix fitting error as defined by a sum of squared euclidean distances from each observed column vector to the subspace spanned by the corresponding rows in  $\mathbf{M}$ .

The error function (2) is minimized using our Column Space Fitting method, based on Levenberg-Marquardt optimization [4], which we summarize in Algorithm 1.  $\mathbf{M}$  is first set to an initial matrix. Then, in each iteration, we update the current estimate of  $\mathbf{M}$  by computing an adjustment matrix  $\Delta \mathbf{M}$  in vectorized form,  $\text{vec}(\Delta \mathbf{M})$ , which stacks the columns of  $\Delta \mathbf{M}$  in a single vector of unknowns. We solve for  $\text{vec}(\Delta \mathbf{M})$  using the gradient vector ( $\mathbf{g}$ ) and Hessian matrix ( $\mathbf{H}$ ) of  $f$  as given by matrix differential calculus [20]. The damping parameter  $\delta$  leads to combined Gauss-Newton and steepest-descent iterations when  $\mathbf{H}$

becomes singular. For numerical stability, we orthogonalize  $\mathbf{M}$  at the end of each iteration.

**Algorithm 1.** Column Space Fitting (CSF) method for minimizing the error function  $f(\mathbf{M})$ .

- 1:  $\mathbf{M} \leftarrow$  initial matrix ( $\mathbf{M}_0$ ).
- 2:  $\delta \leftarrow$  initial damping scalar ( $\delta_0$ ).
- 3: **repeat**
- 4: Compute gradient ( $\mathbf{g}$ ) and Hessian ( $\mathbf{H}$ ) from Jacobian terms ( $\mathbf{J}_j$ ).
- 5: **repeat**
- 6:  $\delta \leftarrow \delta \times 10$ .
- 7: Find  $\Delta \mathbf{M}$  from  $\text{vec}(\Delta \mathbf{M}) \leftarrow (\mathbf{H} + \delta \mathbf{I})^{-1} \mathbf{g}$ .
- 8: **until**  $f(\mathbf{M} - \Delta \mathbf{M}) < f(\mathbf{M})$ .
- 9:  $\mathbf{M} \leftarrow \mathbf{M} - \Delta \mathbf{M}$ .
- 10:  $\delta \leftarrow \delta \times 10^{-2}$ .
- 11: Orthogonalize  $\mathbf{M}$  (keep mean vector  $\mathbf{t}$  unchanged).
- 12: **until** convergence.

To compute  $\mathbf{g}$  and  $\mathbf{H}$ , we follow a Gauss-Newton-based derivation in terms of Jacobian matrices  $\mathbf{J}_j$  for each vector  $\mathbf{w}_j$  (see the complete derivation in Appendix A, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>). The first and second differentials of  $f(\mathbf{M})$  are

$$df = - \sum_j (\mathbf{J}_j^T \mathbf{r}_j)^T \text{vec}(d\mathbf{M}), \quad (5)$$

$$d^2 f = \sum_j \text{vec}(d\mathbf{M})^T (\mathbf{J}_j^T \mathbf{J}_j) \text{vec}(d\mathbf{M}), \quad (6)$$

$$\mathbf{J}_j = \mathbf{s}_j^T \otimes \mathbf{P}_j^\perp \mathbf{\Pi}_j. \quad (7)$$

In (7),  $\otimes$  is the Kronecker product and  $\mathbf{s}_j = \mathbf{M}_j^\dagger \mathbf{w}_j$  denotes the current (and implicit) estimate of the  $j$ th column of factor  $\mathbf{S}$ . From (5) and (6), we identify

$$\mathbf{g} = - \sum_j \mathbf{J}_j^T \mathbf{r}_j \quad \text{and} \quad \mathbf{H} = \sum_j \mathbf{J}_j^T \mathbf{J}_j. \quad (8)$$

For comparison, the forms of  $\mathbf{g}$  and  $\mathbf{H}$  as computed in [8] are reproduced in our supplementary documentation file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>. It is important to note that the simplicity of our expressions above makes it easier to consider additional optimization constraints as described next.

As an alternative to the derivations given above, we can solve for rank- $(r-1)$  factors  $\mathbf{M}$  and  $\mathbf{S}$  and an additional mean column  $\mathbf{t} \in \mathbb{R}^m$  such that  $\mathbf{W} = \mathbf{M} \mathbf{S} + \mathbf{t} \mathbf{1}^T$ , where  $\mathbf{1} \in \mathbb{R}^n$  is a vector of all ones. This is a very useful model in applications such as rigid and nonrigid SFM. Equivalently, we solve for a rank- $r$  factorization

$$\mathbf{W} = \widetilde{\mathbf{M}} \widetilde{\mathbf{S}} = [\mathbf{M} \ \mathbf{t}] \begin{bmatrix} \mathbf{S} \\ \mathbf{1}^T \end{bmatrix}, \quad (9)$$

with extended factors  $\widetilde{\mathbf{M}}$  and  $\widetilde{\mathbf{S}}$ . The mean column interpretation for the last column of  $\widetilde{\mathbf{M}}$  constrains the last row of  $\widetilde{\mathbf{S}}$  to be  $\mathbf{1}^T$ . Therefore, the model in (9) has fewer

degrees of freedom than that in (1) and the final matrix fitting error,  $f(\tilde{\mathbf{M}})$ , is expected to be higher.

The goal now is to minimize

$$f(\mathbf{M}, \mathbf{t}) = \frac{1}{2} \sum_j \mathbf{r}_j^T \mathbf{r}_j, \quad \mathbf{r}_j = \mathbf{P}_j^\perp (\mathbf{w}_j - \mathbf{t}_j), \quad (10)$$

with  $\mathbf{t}_j = \mathbf{\Pi}_j \mathbf{t}$  and  $\mathbf{P}_j^\perp$  computed from  $\mathbf{M}_j$  only. The new Jacobian terms are then derived for  $\text{vec}(d\tilde{\mathbf{M}})$  as

$$\mathbf{J}_j = [\mathbf{s}_j^T \ 1] \otimes \mathbf{P}_j^\perp \mathbf{\Pi}_j, \quad (11)$$

with  $\mathbf{s}_j = \mathbf{M}_j^\dagger (\mathbf{w}_j - \mathbf{t}_j)$  (see the derivation in Appendix A, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>).

In this case, at the end of each iteration, only the leftmost  $r-1$  columns of  $\tilde{\mathbf{M}}$  are orthogonalized and  $\mathbf{t}$  is kept unchanged.

## 2.2 Solution with a Predefined Basis

Let the columns of a matrix  $\mathbf{B} \in \mathbb{R}^{m \times d}$  contain a predefined set of basis vectors (e.g., a truncated DCT basis). We now consider a factor  $\mathbf{M}$  of the form

$$\mathbf{M} = \mathbf{B}\mathbf{X}, \quad \mathbf{X} \in \mathbb{R}^{d \times r}.$$

Each column of  $\mathbf{X}$  has  $d \leq m$  coordinates for the corresponding column of  $\mathbf{M}$  as represented in the  $d$ -dimensional space spanned by  $\mathbf{B}$ .

Previously, we had implicitly considered a canonical basis  $\mathbf{B} = \mathbf{I}_m$ , the  $m \times m$  identity matrix, and solved for  $\mathbf{X} = \mathbf{M}$ . Now, in each iteration we solve for the update step  $\text{vec}(d\mathbf{X})$  with gradient and Hessian as in (8). The Jacobian terms are given by

$$\mathbf{J}_j = \mathbf{s}_j^T \otimes \mathbf{P}_j^\perp \mathbf{\Pi}_j \mathbf{B}. \quad (12)$$

It is important to describe our motivation for using a predefined basis: Problems in different application domains involve matrices ( $\mathbf{W}$ ) whose columns have observations of random variables that change only gradually over time. Observations in each column can be considered as samples of a smooth signal over time with a narrow bandwidth spectrum in the DCT domain. This means that most of the energy of a signal (column) is captured by a small number of low-frequency DCT components, leading to a compact representation  $\mathbf{X}$  (i.e., small  $d$ ).

In the next sections, we consider matrices  $\mathbf{B}$  constructed from DCT basis vectors, with the rightmost columns corresponding to higher frequency components. Since signal smoothness leads to high-frequency coefficients close to zero, we propose an initialization of  $\mathbf{X}$ , denoted  $\mathbf{X}_0$ , of the form

$$\mathbf{X}_0 = \begin{bmatrix} \mathbf{Q} \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{Q} \in \mathbb{R}^{r \times r}.$$

The goal becomes initializing a small block  $\mathbf{Q}$ , which must be full rank (i.e., the rank of  $\mathbf{M} = \mathbf{B}\mathbf{X}_0$  is  $\text{rank}(\mathbf{Q})$ ). For any full-rank  $\mathbf{Q}$ , an equally good initial solution is given by the simple and deterministic initialization

$$\mathbf{X}_0 = \begin{bmatrix} \mathbf{Q} \\ \mathbf{0} \end{bmatrix} \mathbf{Q}^{-1} = \begin{bmatrix} \mathbf{I}_r \\ \mathbf{0} \end{bmatrix}. \quad (13)$$

As a proof, we note that the factorization in (1) is defined only up to a rank- $r$  ambiguity matrix  $\mathbf{Q}$  such that  $\mathbf{W} = \mathbf{M}\mathbf{S} = (\mathbf{M}\mathbf{Q}^{-1})(\mathbf{Q}\mathbf{S})$ .

As a result, we can even fix the values for the topmost block of  $\mathbf{X}$  as  $\mathbf{I}_r$  and solve only for the remaining  $(d-r) \times r$  submatrix. Note that  $(d-r)r$  is the number of parameters defining an  $r$ -dimensional linear space embedded in another  $d$ -dimensional linear space.

When the last column of  $\mathbf{X}$  is interpreted as the coefficients of the mean column vector  $\mathbf{t}$ , the ambiguity matrix  $\mathbf{Q}$  is of rank  $(r-1)$  only, i.e., the last row of  $\mathbf{S}$  is constrained to be  $\mathbf{1}^T$ . This fact suggests initializing the last column of  $\mathbf{X}$  with zeros ( $\mathbf{t} = \mathbf{0}$ ).

## 3 RIGID SFM: ESTIMATING THE SMOOTH TIME TRAJECTORY OF A CAMERA

This section focuses on rigid SFM with occlusion. We assume the point tracks in  $\mathbf{W}$  were obtained from a monocular video sequence provided by a single, smoothly moving camera. First, we present a method using a compact DCT basis to represent the smooth trajectory of a general affine camera. Subsequently, we further improve on this method by using a weak-perspective camera model from the outset to directly estimate euclidean cameras and shape. Once the final solution for  $\mathbf{M}$  is defined, each column  $\mathbf{s}_j$  of  $\mathbf{S}$  is obtained independently as in the computation of the Jacobian terms above.

### 3.1 Affine Camera Trajectory

In rigid SFM with  $T$  images,  $\mathbf{W} \in \mathbb{R}^{2T \times n}$  is

$$\mathbf{W} = \begin{bmatrix} x_{11} & \dots & x_{1n} \\ y_{11} & \dots & y_{1n} \\ \vdots & & \vdots \\ x_{T1} & \dots & x_{Tn} \\ y_{T1} & \dots & y_{Tn} \end{bmatrix} = \mathbf{M} \begin{bmatrix} \mathbf{s}_1 & \dots & \mathbf{s}_n \\ 1 & & 1 \end{bmatrix}, \quad (14)$$

with  $[x_{tj}, y_{tj}]^T$  the 2D projection of 3D point  $\mathbf{s}_j$  in the  $t$ th image ( $t = 1, 2, \dots, T$ ;  $j = 1, 2, \dots, n$ ). Here, the motion factor  $\mathbf{M} \in \mathbb{R}^{2T \times 4}$  describes a stack of  $T$  affine camera matrices  $\widehat{\mathbf{M}}_t \in \mathbb{R}^{2 \times 4}$ ,

$$\mathbf{M} = \begin{bmatrix} \widehat{\mathbf{M}}_1 \\ \vdots \\ \widehat{\mathbf{M}}_T \end{bmatrix}, \quad \widehat{\mathbf{M}}_t = [\widehat{\mathbf{A}}_t \ \widehat{\mathbf{t}}_t], \quad (15)$$

where  $\widehat{\mathbf{A}}_t \in \mathbb{R}^{2 \times 3}$  is a general affine projection and  $\widehat{\mathbf{t}}_t \in \mathbb{R}^2$  is a 2D translation. The rows of  $\widehat{\mathbf{A}}_t$  are 3D vectors defining the  $x$  and  $y$ -axes (not necessarily orthonormal) of the camera's projection plane in the  $t$ th observation.

Assuming the eight camera parameters vary smoothly over time, we define the camera as a sample of a smooth matrix function of time,  $\widehat{\mathbf{M}}_t = \widehat{\mathbf{M}}(t)$ . We then parameterize  $\widehat{\mathbf{M}}(t)$  using eight independent cosine series whose  $f$ th frequency coefficients are given in  $\widehat{\mathbf{X}}_f \in \mathbb{R}^{2 \times 4}$ :

$$\widehat{\mathbf{M}}(t) = [\widehat{\mathbf{A}}(t) \ \widehat{\mathbf{t}}(t)] = \sum_{f=1}^d \omega_{tf} \widehat{\mathbf{X}}_f. \quad (16)$$

Each constant  $\omega_{tf}$  above is the  $f$ th frequency cosine term at time  $t$ :

$$\omega_{tf} = \frac{\sigma_f}{\sqrt{T}} \cos\left(\frac{\pi(2t-1)(f-1)}{2T}\right),$$

with  $\sigma_1 = 1$  and, for  $f \geq 2$ ,  $\sigma_f = \sqrt{2}$ .

The goal is then to solve for  $d \leq T$  matrices  $\widehat{\mathbf{X}}_f$ . Let  $\Omega_d \in \mathbb{R}^{T \times d}$  be a truncated, orthonormal DCT matrix whose  $(t, f)$ th entry is  $\omega_{tf}$ , as above. Substituting (16) into (15) yields

$$\mathbf{M} = (\Omega_d \otimes \mathbf{I}_2) \begin{bmatrix} \widehat{\mathbf{X}}_1 \\ \vdots \\ \widehat{\mathbf{X}}_d \end{bmatrix} = \mathbf{B}_{df} \mathbf{X}.$$

Thus, with basis  $\mathbf{B}_{df} = (\Omega_d \otimes \mathbf{I}_2) \in \mathbb{R}^{2T \times 2d}$ , we solve for  $\mathbf{M} \in \mathbb{R}^{2T \times 4}$  as a function of  $\mathbf{X} \in \mathbb{R}^{2d \times 4}$ . Note that the smoothness assumption implies  $d \ll T$ . We solve for  $\mathbf{X}$  as in Section 2.2, with extended terms  $\mathbf{s}_j^T = [\mathbf{s}_j^T \ 1]$  in (12).

For simplicity, and unless stated otherwise, we assume in this paper that the number of DCT components  $d$  has been fixed a priori based on empirical observations (e.g., expected motion or noise level). An alternative *coarse-to-fine* strategy can start with a small  $d$  and increment its value until a convergence criterion is verified. For instance,  $d$  can be incremented until the energy of the highest frequency coefficients falls below a small percentage  $p_d \in (0, 1)$  of the total energy,

$$\|\mathbf{x}_d^T\|_F < p_d \|\mathbf{X}\|_F,$$

where  $\mathbf{x}_d^T$  is the last row of  $\mathbf{X}$ .

### 3.2 Euclidean Camera Trajectory

The solution  $\mathbf{X}$  above leads to an affine shape  $\mathbf{S}$ . Recovering euclidean shape and cameras requires a subsequent upgrading step and, ideally, a final bundle adjustment [15]. Because an affine camera has more degrees of freedom than the final euclidean camera, the initial solution can overfit measurement noise and also lead to incorrect results in cases of degenerate, planar camera motion [21]. These are problems faced by most SFM methods based on an initial affine solution. To avoid them and also render postprocessing unnecessary, we further improve on our method above.

We now solve directly for the smooth trajectory of a euclidean camera, also obtaining the optimal euclidean shape directly. We consider the trajectory of a weak-perspective (i.e., scaled orthographic) camera:

$$\widehat{\mathbf{M}}_t = \begin{bmatrix} \widehat{\mathbf{R}}_t & \widehat{\mathbf{t}}_t \end{bmatrix},$$

where the vector  $\widehat{\mathbf{t}}_t$  is defined as before. The rows of  $\widehat{\mathbf{R}}_t \in \mathbb{R}^{2 \times 3}$  define the euclidean 2D imaging plane of the camera at time  $t$ . These rows are constrained to be two equal length, orthogonal 3D vectors and are obtained from a scaled 3D rotation:

$$\widehat{\mathbf{R}}_t = \lambda_t \begin{bmatrix} \cos \alpha_t \cos \beta_t \cos \gamma_t - \sin \alpha_t \sin \gamma_t & & \\ -\cos \alpha_t \cos \beta_t \sin \gamma_t - \sin \alpha_t \cos \gamma_t & \cos \alpha_t \sin \beta_t & \\ \sin \alpha_t \cos \beta_t \cos \gamma_t + \cos \alpha_t \sin \gamma_t & & \\ -\sin \alpha_t \cos \beta_t \sin \gamma_t + \cos \alpha_t \cos \gamma_t & \sin \alpha_t \sin \beta_t & \end{bmatrix}. \tag{17}$$

Hence,

$$\widehat{\mathbf{R}}_t = \lambda_t \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \mathbf{R}_Z(\alpha_t) \mathbf{R}_Y(\beta_t) \mathbf{R}_Z(\gamma_t).$$

The scalar  $\lambda_t$  is the weak-perspective scale and the three Euler angles  $\alpha_t$ ,  $\beta_t$ , and  $\gamma_t$  determine a 3D rotation as a sequence of simpler rotations around the  $Z$  and  $Y$  world coordinate axes.

We model a smooth camera trajectory by considering a cosine series for each of the six camera parameters given by the weak-perspective model above,

$$\begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_T \end{bmatrix} = \Omega_d \mathbf{x}_1, \quad \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_T \end{bmatrix} = \Omega_d \mathbf{x}_2, \quad \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_T \end{bmatrix} = \Omega_d \mathbf{x}_3,$$

$$\begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_T \end{bmatrix} = \Omega_d \mathbf{x}_4, \quad \mathbf{t} = \mathbf{B}_{df} \mathbf{x}_5,$$

where vectors  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4 \in \mathbb{R}^d$  and  $\mathbf{x}_5 \in \mathbb{R}^{2d}$  have the unknown DCT coefficients. For simplicity of presentation, the notation above assumes that  $d$  is the same for  $\mathbf{x}_1, \dots, \mathbf{x}_5$ . Note, however, that our approach is not limited to such a case—an alternative derivation is given in the supplementary file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50> for the case in which  $\mathbf{x}_1, \dots, \mathbf{x}_5$  have different numbers of DCT components.

A deterministic initial solution is defined with  $\mathbf{t} = \mathbf{0}$  ( $\mathbf{x}_5 = \mathbf{0}$ ), as before. Also note that the rigid SFM solution we seek here is defined only up to an arbitrary scale and a  $3 \times 3$  rotation. We thus initialize the sequence of scaling factors  $\lambda_t$  to a constant signal by setting  $\mathbf{x}_4 = [1, 0, \dots, 0]^T$ . We avoid initializing the angles  $\beta_t$  ( $\mathbf{x}_2$ ) with zeros because the resulting  $\widehat{\mathbf{R}}_t$  would then represent only 2D rotations within a plane. Setting  $\mathbf{x}_2 = [0, 1, 0, \dots, 0]^T$  provides a coarse initial solution with a smooth sequence of *noncoincident* camera planes that allows  $\mathbf{x}_1$  and  $\mathbf{x}_3$  to be initialized with zeros.

Here, we iteratively solve for  $\text{vec}(d\mathbf{X})$  defined as

$$\text{vec}(d\mathbf{X}) \equiv [d\mathbf{x}_1^T, d\mathbf{x}_2^T, d\mathbf{x}_3^T, d\mathbf{x}_4^T, d\mathbf{x}_5^T]^T \in \mathbb{R}^{6d}.$$

To this end, modified expressions for  $\text{vec}(d\mathbf{M})$  and the Jacobian terms are needed. First, let  $\omega_t^T$  be the  $t$ th row of  $\Omega_d$ . The differential of each camera parameter at time  $t$  is  $d\alpha_t = \omega_t^T d\mathbf{x}_1$ ,  $d\beta_t = \omega_t^T d\mathbf{x}_2$ ,  $d\gamma_t = \omega_t^T d\mathbf{x}_3$ ,  $d\lambda_t = \omega_t^T d\mathbf{x}_4$ , and  $d\mathbf{t} = (\omega_t^T \otimes \mathbf{I}_2) d\mathbf{x}_5$ . The differential of each rotation matrix (17) is

$$d\widehat{\mathbf{R}}_t = \frac{\partial \widehat{\mathbf{R}}_t}{\partial \alpha_t} d\alpha_t + \frac{\partial \widehat{\mathbf{R}}_t}{\partial \beta_t} d\beta_t + \frac{\partial \widehat{\mathbf{R}}_t}{\partial \gamma_t} d\gamma_t + \frac{\partial \widehat{\mathbf{R}}_t}{\partial \lambda_t} d\lambda_t.$$

Considering the three columns  $\widehat{\mathbf{r}}_{t,1}$ ,  $\widehat{\mathbf{r}}_{t,2}$ , and  $\widehat{\mathbf{r}}_{t,3} \in \mathbb{R}^2$  of a partial derivative matrix  $\frac{\partial \widehat{\mathbf{R}}_t}{\partial \alpha_t}$ , for all  $t$ , we stack all  $\alpha$ -terms associated with  $d\mathbf{x}_1$  into

$$\mathbf{B}_\alpha = \begin{bmatrix} \hat{\mathbf{r}}_{1,1} \otimes \boldsymbol{\omega}_1^T \\ \vdots \\ \hat{\mathbf{r}}_{T,1} \otimes \boldsymbol{\omega}_T^T \\ \hat{\mathbf{r}}_{1,2} \otimes \boldsymbol{\omega}_1^T \\ \vdots \\ \hat{\mathbf{r}}_{T,2} \otimes \boldsymbol{\omega}_T^T \\ \hat{\mathbf{r}}_{1,3} \otimes \boldsymbol{\omega}_1^T \\ \vdots \\ \hat{\mathbf{r}}_{T,3} \otimes \boldsymbol{\omega}_T^T \end{bmatrix} \in \mathbb{R}^{6T \times d}.$$

Analogously, we define  $\mathbf{B}_\beta$ ,  $\mathbf{B}_\gamma$ , and  $\mathbf{B}_\lambda$ . Hence,

$$\begin{aligned} \text{vec}(d\mathbf{M}) &= \underbrace{\begin{bmatrix} \mathbf{B}_\alpha & \mathbf{B}_\beta & \mathbf{B}_\gamma & \mathbf{B}_\lambda & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{B}_{af} \end{bmatrix}}_{\tilde{\mathbf{B}}_{wp}} \begin{bmatrix} dx_1 \\ dx_2 \\ dx_3 \\ dx_4 \\ dx_5 \end{bmatrix} \\ &= \tilde{\mathbf{B}}_{wp} \text{vec}(d\mathbf{X}), \end{aligned}$$

with  $\tilde{\mathbf{B}}_{wp} \in \mathbb{R}^{8T \times 6d}$  and  $\mathbf{0}$  is a matrix of zeros.

Finally, from (5) to (7), the Jacobian terms of this camera model can be computed as

$$\mathbf{J}_j = ([s_j^T \ 1] \otimes \mathbf{P}_j^\perp \boldsymbol{\Pi}_j) \tilde{\mathbf{B}}_{wp}.$$

Unlike in the previous sections, here the extended basis  $\tilde{\mathbf{B}}_{wp}$  must be recomputed in each iteration to update the partial derivatives in  $\mathbf{B}_\alpha$ ,  $\mathbf{B}_\beta$ ,  $\mathbf{B}_\gamma$ , and  $\mathbf{B}_\lambda$ . Thus,  $\tilde{\mathbf{B}}_{wp}$  can be considered only as a local basis at the current location of the smooth parameter manifold.

The last step in each iteration shown in Algorithm 1, orthogonalization of  $\mathbf{M}$ , is no longer necessary.

#### 4 NONRIGID SFM: ESTIMATING A SMOOTH TRAJECTORY IN SHAPE SPACE

In this section, we solve for nonrigid SFM by estimating the smooth time trajectory of a 3D shape as represented by a point moving in a linear shape space. Without loss of generality, we assume  $\mathbf{W}$  is complete and derive our 3D shape trajectory approach. Subsequently, we offer an algorithm for cases with missing data.

##### 4.1 Solving for a Shape Trajectory

Let  $\mathbf{W} \in \mathbb{R}^{2T \times n}$  as in (14) and consider the rank- $r$  factorization method, with  $r = 3K + 1$ , proposed by Bregler et al. [5] for nonrigid scenes,

$$\mathbf{W} = \underbrace{\mathbf{D}(\mathbf{C} \otimes \mathbf{I}_3)}_{\mathbf{M}} \mathbf{S} + \mathbf{t}\mathbf{1}^T,$$

with  $\mathbf{D} \in \mathbb{R}^{2T \times 3T}$ ,  $\mathbf{C} \in \mathbb{R}^{T \times K}$ , and  $\mathbf{S} \in \mathbb{R}^{3K \times n}$ . The predefined constant  $K$  is the number of 3D basis shapes  $\hat{\mathbf{S}}_k \in \mathbb{R}^{3 \times n}$  defining a linear shape space in  $\mathbf{S} = [\hat{\mathbf{S}}_1^T \ \hat{\mathbf{S}}_2^T \ \dots \ \hat{\mathbf{S}}_K^T]^T$ .

The column space factor  $\mathbf{M}$  is composed of a block-diagonal rotation matrix  $\mathbf{D}$ :

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{R}}_1 & & & \\ & \hat{\mathbf{R}}_2 & & \\ & & \ddots & \\ & & & \hat{\mathbf{R}}_T \end{bmatrix},$$

and a shape coordinate matrix  $\mathbf{C}$ . Each row  $\mathbf{c}_t^T \in \mathbb{R}^K$  of  $\mathbf{C}$  has the coordinates of the 3D shape in the  $t$ th image with respect to the shape basis in  $\mathbf{S}$ . Here, we also consider  $\mathbf{c}_t^T = c(t)$  as a single point in shape space that defines a *single smooth 3D shape trajectory over time*.

We assume that the smooth 2D trajectories in  $\mathbf{W}$  reflect not only smooth camera motion, but also the smooth deformation of the observed 3D structure over time. This means that each shape coordinate  $c_{tk} \in \mathbb{R}$  ( $k = 1, 2, \dots, K$ ) is assumed to vary smoothly with  $t$ . Then, we represent  $\mathbf{C}$  using  $K$  compact cosine series:

$$\mathbf{C} = \begin{bmatrix} c_{1,1} & \dots & c_{1,K} \\ \vdots & \ddots & \vdots \\ c_{T,1} & \dots & c_{T,K} \end{bmatrix} = \boldsymbol{\Omega}_d \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_K \end{bmatrix},$$

with  $\mathbf{x}_k \in \mathbb{R}^d$  ( $k = 1, \dots, K$ ),

$$\mathbf{C} = \boldsymbol{\Omega}_d \mathbf{X}, \quad \mathbf{X} \in \mathbb{R}^{d \times K}.$$

Assuming  $\mathbf{W}$  is complete,  $\mathbf{t}$  is estimated simply as the mean column of  $\mathbf{W}$ . Then, to recover euclidean shapes and cameras, we consider  $\mathbf{S}$  implicitly and solve for

$$\mathbf{M} = \mathbf{D}(\boldsymbol{\Omega}_d \mathbf{X} \otimes \mathbf{I}_3), \quad (18)$$

with  $\mathbf{D}$  subject to camera orthonormality constraints, i.e.,  $\hat{\mathbf{R}}_t \hat{\mathbf{R}}_t^T = \mathbf{I}_2, \forall t$ . For now, let's assume  $\mathbf{D}$  has been computed by an initialization algorithm. We will define this algorithm below. Thus, we only need to solve for the rank- $K$  shape trajectory  $\mathbf{X}$  in the DCT domain.

For a fixed  $\mathbf{D}$ , the factor  $\mathbf{C} = \boldsymbol{\Omega}_d \mathbf{X}$  is defined only up to a full-rank ambiguity  $\mathbf{Q} \in \mathbb{R}^{K \times K}$ ; equivalently,  $\mathbf{M}$  in (18) is defined only up to an ambiguity  $\mathbf{Q} \otimes \mathbf{I}_3$ . Therefore, as in Section 2.2, we can initialize  $\mathbf{X}$  with a coarse solution  $\mathbf{X}_0 = [\mathbf{I}_K \ \mathbf{0}]^T$ , leading to

$$\mathbf{M}_0 = \mathbf{D}(\boldsymbol{\Omega}_d \mathbf{X}_0 \otimes \mathbf{I}_3) = \mathbf{D}(\boldsymbol{\Omega}_K \otimes \mathbf{I}_3). \quad (19)$$

Note that  $K < d$  and the initial rank- $3K$  solution in (19) can only use  $K$  low-frequency vectors in the DCT basis.

The 3D point trajectory approach (PTA) in [2] defines  $\mathbf{M} = \mathbf{D}\boldsymbol{\Theta}$ , where  $\boldsymbol{\Theta}$  has the same columns as  $(\boldsymbol{\Omega}_K \otimes \mathbf{I}_3)$  in (19), but in a different order. Therefore, our coarse initial solution  $\mathbf{M}_0$  is equivalent to the final solution  $\mathbf{M}$  of PTA (with factor  $\mathbf{S}$  presenting a different order of rows). Note that the PTA method cannot consider additional DCT basis vectors without increasing  $K$ , leading to a higher rank of  $\mathbf{M}$ . Our 3D shape trajectory approach, on the other hand, can consider any number  $d = K, \dots, T$  of DCT basis vectors because the linear combination represented by  $\mathbf{X}$  constrains  $\mathbf{M}$  to be of rank- $3K$ . That is, our method can better model structure deformation presenting higher frequency components in the DCT domain, yielding better 3D shape reconstructions.

Empirically, we have found that the coarse solutions computed by PTA contain accurate estimates of the rotation matrices in  $\mathbf{D}$ . We iteratively run PTA with increasing

values of  $K \in \{1, 2, \dots, \lfloor \frac{n}{3} \rfloor\}$ , obtaining a solution denoted as  $\mathbf{D}_K$ . Iterations stop automatically when there is no additional improvement in the average camera orthonormality:

$$\varepsilon(\mathbf{D}_K) = \frac{1}{T} \sum_{t=1}^T \|\mathbf{I}_2 - \widehat{\mathbf{R}}_t \widehat{\mathbf{R}}_t^T\|_F^2.$$

Given  $\mathbf{D} = \mathbf{D}_K$ , we then solve for  $\mathbf{M}$  as a function of  $\mathbf{X}$  only. The first differential of  $\mathbf{M}$  in (18) is given by

$$d\mathbf{M} = \underbrace{\mathbf{D}(\boldsymbol{\Omega}_d \otimes \mathbf{I}_3)}_{\mathbf{B}_{nr}}(d\mathbf{X} \otimes \mathbf{I}_3) = \mathbf{B}_{nr}(d\mathbf{X} \otimes \mathbf{I}_3). \quad (20)$$

Next, consider  $\text{vec}(d\mathbf{X} \otimes \mathbf{I}_3) = \mathbf{V} \text{vec}(d\mathbf{X})$ , with  $\mathbf{V} \in \mathbb{R}^{9dK \times dK}$  a binary mapping matrix. From (12) and (20), the Jacobian terms for the update step  $\text{vec}(d\mathbf{X})$  are

$$\mathbf{J}_j = (\mathbf{s}_j^T \otimes \mathbf{P}_j^+ \boldsymbol{\Pi}_j \mathbf{B}_{nr}) \mathbf{V},$$

with  $\mathbf{s}_j = \mathbf{M}_j^\dagger(\mathbf{w}_j - \mathbf{t}_j)$ . For  $\mathbf{X}$  with dimensions  $d \times K$ , the constant and sparse matrix  $\mathbf{V}$  is defined as [20]

$$\mathbf{V} = \mathbf{I}_K \otimes [(\mathbf{K}_{3d} \otimes \mathbf{I}_3)(\mathbf{I}_d \otimes \text{vec}(\mathbf{I}_3))],$$

where the permutation matrix  $\mathbf{K}_{3d} \in \mathbb{R}^{3d \times 3d}$  satisfies  $\text{vec}(\mathbf{A}^T) = \mathbf{K}_{3d} \text{vec}(\mathbf{A})$ , for any  $\mathbf{A} \in \mathbb{R}^{3 \times d}$ .

### 4.2 Nonrigid SFM with Missing Data

The method above can estimate  $\mathbf{X}$  using only the observed data in  $\mathbf{W}$ . We now consider the initial estimation of  $\mathbf{t}$  and  $\mathbf{D}$  (the camera motion) in cases of missing data. To compute  $\mathbf{t}$  and  $\mathbf{D}$  from a coarse solution as above, we first recover a complete, rank- $r$   $\mathbf{W}$  via the factorization

$$\mathbf{W} = \underbrace{\mathbf{B}_{af} \mathbf{X}}_{\mathbf{M}} \begin{bmatrix} \mathbf{S} \\ \mathbf{1}^T \end{bmatrix}, \quad \mathbf{X} \in \mathbb{R}^{2d \times r},$$

with a predefined  $r \in \{4, 5, \dots, 3K + 1\}$  and the DCT basis  $\mathbf{B}_{af}$  of Section 3. Instead of considering the trajectory of an affine camera, here  $\mathbf{B}_{af}$  is a basis for *individual, smooth 2D point trajectories* in the column space of  $\mathbf{W}$ .

Note that a common solution for  $\mathbf{M}$  corresponds to the eigenvector matrix  $\mathbf{U}$  of  $\mathbf{W}\mathbf{W}^T$  (if  $\mathbf{W}$  is complete). Here, our solution  $\mathbf{M} = \mathbf{B}_{af} \mathbf{X}$  can be seen as using linear combinations of the DCT basis vectors to approximate the eigenvectors in  $\mathbf{U}$ . Indeed, the DCT basis vectors have been found to be good approximations of the Karhunen-Loeve (eigenfunction) transform of first-order Markov processes [18]. Also, the Markov assumption is used in algorithms that track the 2D points in  $\mathbf{W}$  [11].

## 5 EXPERIMENTAL RESULTS

In this section, we provide extensive experimental validation for the proposed algorithms. General matrix factorization performance is first assessed on different, synthetic and real data sets. Subsequently, we apply our methods on rigid and nonrigid SFM data.

### 5.1 Fitting Low-Rank Matrices

Our initial experiments analyze the general performance of our CSF method in its simplest form (i.e., with a canonical basis  $\mathbf{I}$  and without a mean column vector  $\mathbf{t}$ ) in the low-rank

TABLE 1  
Algorithm Performances on Random Matrices with Missing Data: Frequency of Suboptimal Solutions (Percent) and Average Number of Iterations (in Parentheses)

$\rho(\%), \sigma_n$	PF	LM-S	LM-S <sub>GN</sub>	CSF
25, .1	23 (313)	<b>1 (16)</b>	4 (82)	3 (57)
25, .2	30 (391)	<b>3 (17)</b>	6 (93)	<b>3 (65)</b>
25, .4	28 (381)	<b>3 (16)</b>	8 (108)	5 (74)
50, .1	71 (813)	12 (51)	14 (96)	<b>11 (79)</b>
50, .2	76 (831)	15 (58)	14 (105)	<b>12 (89)</b>
50, .4	77 (869)	18 (61)	19 (120)	<b>16 (103)</b>
75, .1	100 (1000)	69 (356)	62 (380)	<b>42 (549)</b>
75, .2	100 (1000)	72 (376)	54 (419)	<b>46 (586)</b>
75, .4	100 (1000)	77 (421)	65 (473)	<b>44 (698)</b>

factorization of synthetic and real data matrices with missing data. The method is evaluated against the LM-S algorithm [8] and its Gauss-Newton variant, LM-S<sub>GN</sub> (presented in the supplementary file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>). We also provide results using PowerFactorization [16] as a baseline algorithm.

We start by generating a random  $20 \times 30$  matrix  $\mathbf{W}$  of rank 3 with values uniformly distributed in the interval  $[0, 1]$ . Then, we add Gaussian noise with standard deviation  $\sigma_n$  and randomly occlude  $\rho$  percent of the matrix entries. The algorithms above are used to compute a rank-3 factorization  $\mathbf{W} = \mathbf{M}\mathbf{S}$  starting from the same initial factor  $\mathbf{M} = \mathbf{M}_0$ , which is generated randomly and then refined with 20 iterations of the PF method. We always set the initial damping scalar to  $\delta_0 = 10^{-4}$ . Each algorithm runs until the change in the cost value (2) is less than  $10^{-10}$  or until the number of iterations reaches 1,000. Note that, although we know the ground truth data, we do not know where the optimal solution for the matrix with missing data is [8]. Therefore, we first run all algorithms and, considering the solution obtained by each algorithm, we then define the trial’s target (i.e., “optimal”) cost as the smallest cost observed.

We perform 500 trials, with  $\sigma_n$  and  $\rho$  fixed, and report the frequency (percent) with which a method failed to match the trial’s “optimal cost.” We consider two cost values as equivalent if their absolute difference is below  $10^{-7}$ . Results for different values of  $\sigma_n$  and  $\rho$  are shown in Table 1, which also includes the average number of iterations performed by each method. As expected, the simple PF method converges very slowly and is competitive only when  $\rho$  is very small. For  $\rho \in \{25\%, 50\%\}$ , CSF’s performance is either comparable or better than that of LM-S and LM-S<sub>GN</sub>. CSF also converges faster than LM-S<sub>GN</sub> for smaller fractions of missing data. With  $\rho = 75\%$ , CSF clearly outperforms the other methods. These results are highly relevant to SFM, where the number of missing entries is usually large ( $> 75\%$ ).

A second, similar experiment considers the three real data sets used in [6] and [8]: 1) the *dinosaur* on a turn table, with rank-4  $\mathbf{W} \in \mathbb{R}^{72 \times 319}$  and 76.9 percent missing data, 2) the occluded motion of a *giraffe*, with rank-6  $\mathbf{W} \in \mathbb{R}^{240 \times 166}$  and 30.2 percent missing data, and 3) the *face* illumination data under a moving light source, with rank-4

TABLE 2

Algorithm Performances on Real Data Sets with Missing Data: Frequency of Suboptimal Convergence (Percent) and Average Number of Iterations (in Parenthesis)

Dataset	PF	LM-S	LM- $S_{GN}$	CSF	CSF- $B_{af}$
Dinosaur	100 (500)	26 (225)	3 (92)	<b>1 (101)</b>	<b>0 (12)</b>
Giraffe	100 (500)	0 (192)	0 (86)	<b>0 (61)</b>	<b>0 (49)</b>
Face	100 (500)	65 (36)	62 (64)	<b>62 (43)</b>	-

$\mathbf{W} \in \mathbb{R}^{20 \times 2944}$  and 41.7 percent missing data. For each data set, the smallest known root-mean-square error (RMSE) of the observed entries, as previously published [6], gives the target value that indicates optimal convergence.

We performed 100 trials corresponding to 100 different initial factors  $\mathbf{M}$  generated as above. Each algorithm was allowed to run for at most 500 iterations. Table 2 shows the frequency of suboptimal convergence (percent) and average number of iterations for each algorithm. We first note that on the face data set, all algorithms except PF often converged to a suboptimal solution whose cost is only 0.45 percent (about  $10^{-4}$ ) higher than that of the optimal solution. If this solution is considered equivalent to the optimal one, the new suboptimal convergence frequencies are 9 (LM-S), 8 (LM- $S_{GN}$ ), and 1 percent (CSF). Nevertheless, in the results above, CSF either provides better solutions or computes equivalent solutions with fewer iterations as compared to the other methods.

Because the columns of the dinosaur and giraffe matrices have measurements (alternated  $x$  and  $y$ -coordinates) that vary smoothly over subsequent pairs of rows, we also computed factorizations with the CSF- $B_{af}$  variant of our method. We used the full  $B_{af}$  basis and started with the deterministic initialization in (13) as a coarse initial solution. CSF- $B_{af}$  presented optimal convergence on both data sets, without the initial use of PF.

## 5.2 Computing Rigid SFM

Our first experiment on rigid SFM considers a synthetic *sphere* data set with known 3D points ( $n = 100$ ) located on its surface and also known camera matrices ( $T = 90$ ) describing a smooth trajectory over time, Fig. 3. The real occlusion pattern for the simple spherical shape can be easily determined and the resulting observation matrix  $\mathbf{W} \in \mathbb{R}^{180 \times 100}$  is missing 50 percent of its entries.

Results of our SFM methods with affine (CSF- $B_{af}$ ) and weak-perspective (CSF- $B_{wp}$ ) cameras are compared against those of PF, LM-S, LM- $S_{GN}$ , and the Wiberg algorithm [31]. Euclidean upgrade is applied to the results of affine methods. All results are aligned (rotated and scaled) with

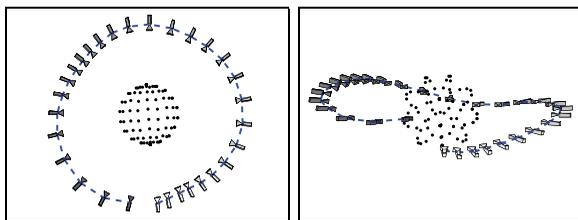


Fig. 3. Smooth camera trajectory in the synthetic sphere data set.

TABLE 3

Results of Rigid SFM Methods on the Sphere Data Set

Method	$\hat{\sigma}_n = .05$		$\hat{\sigma}_n = .125$		$\hat{\sigma}_n = .25$	
	$e_S$	$e_R$	$e_S$	$e_R$	$e_S$	$e_R$
PF	.0163	.0250	.0818	.1247	.2092	.2839
LM-S	.0169	.0255	.0845	.1284	.1780	.2690
LM- $S_{GN}$	.0170	.0255	.0845	.1284	.1777	.2686
Wiberg	.0163	.0250	.0806	.1242	.1695	.2569
CSF- $B_{af}$	.0163	.0151	.0806	<b>.0731</b>	.1651	<b>.1513</b>
CSF- $B_{wp}$	.0156	.0130	.0779	<b>.0606</b>	.1586	<b>.1280</b>

the ground-truth shape and motion before comparison. On the noiseless  $\mathbf{W}$ , the methods above provide camera and shape reconstructions with nearly zero error. Our CSF methods are run with only 30 percent of the DCT basis vectors ( $d = 0.3T$ ).

In our experiment, we analyze how the performance of these methods degrades with different levels of Gaussian noise added to  $\mathbf{W}$ . We now set  $\sigma_n = \hat{\sigma}_n \sigma(\mathbf{W})$ , where  $\hat{\sigma}_n \in [0.05, 0.25]$  and  $\sigma(\mathbf{W})$  indicates the scale of the entries of  $\mathbf{W}$ . For any matrix  $\mathbf{A}$  with  $m$  rows, let

$$\sigma(\mathbf{A}) = \frac{1}{m} \sum_{r=1}^m \sigma_r, \quad (21)$$

where  $\sigma_r$  is the standard deviation of the available entries in the  $r$ th row of  $\mathbf{A}$ . For each value of  $\hat{\sigma}_n$ , we perform 100 trials and report the average error of the recovered 3D shape ( $e_S$ ) and camera rotations ( $e_R$ ). The error  $e_S$  is the average euclidean distance between original and recovered 3D points, normalized by the radius of the original 3D sphere. Let the original and estimated rotations be  $\hat{\mathbf{R}}_t^*$  and  $\hat{\mathbf{R}}_t$ , then  $e_R$  is

$$e_R = \frac{1}{T} \sum_{t=1}^T \left\| \hat{\mathbf{R}}_t^* - \hat{\mathbf{R}}_t \right\|_F. \quad (22)$$

CSF- $B_{af}$  and CSF- $B_{wp}$  were run with deterministic initializations. The other methods require random initializations and often provided very poor solutions in our experiments. For this reason, in each trial they were run with five different random initializations; the result with the smallest RMSE for the reconstructed  $\mathbf{W}$  was chosen.

The average errors in Table 3 show that the shape estimates of all methods seem to be similarly affected by noise up to  $\hat{\sigma}_n = 0.125$ . Also, a significant difference is seen for the camera estimates with  $\hat{\sigma}_n > 0.05$ . Affine methods overfit noise in the data due to the extra degrees of freedom in their camera model. CSF- $B_{af}$  and CSF- $B_{wp}$  attenuate this problem by enforcing smoothness on the camera trajectory. CSF- $B_{wp}$  outperforms all methods by further constraining the camera axes to be orthogonal and of equal length. This result is of special importance to SFM applications that rely more on the motion factor.

The average runtimes in seconds (per initialization, on a single-core 2.6 GHz processor) were: 0.2 (PF), 18.2 (LM-S), 14.9 (LM- $S_{GN}$ ), 128.5 (Wiberg), 1.1 (CSF- $B_{af}$ ), and 8.0 (CSF- $B_{wp}$ ). The simple PF is fast but often provides poor estimates, even with a maximum of 10,000 iterations (versus

TABLE 4  
Results of Rigid SFM Methods on Different Subsets of the Dinosaur Data Set

Method	(model)	319 point tracks (76.9% missing)		2,683 point tracks (87.8% missing)		4,983 point tracks (90.8% missing)	
		RMSE	mean/max 2D error	RMSE	mean/max 2D error	RMSE	mean/max 2D error
CSF- $B_{af}$	(1)	1.0847	.6286 / 38.8381	1.3370	.7129 / 39.9999	1.1346	.5133 / 39.9999
LM-S, LM- $S_{GN}$	(1)	1.0847	.6287 / 38.8376	1.3370	.7131 / 39.9999	1.1346	.5135 / 39.9997
DN	(1)	1.0847	–	–	–	–	–
DP	(1)	1.5482	1.0131 / 39.4277	2.0201	1.9103 / 40.2124	1.7142	1.6233 / 40.1713
MME	(1)	–	–	–	2.4017 / 72.4467	–	1.8438 / 72.4467
CSF- $B_{af}$	(2)	1.2702	.8810 / 41.4078	1.4294	.8610 / 42.5002	1.2176	.6808 / 42.4961
Wiberg	(2)	1.2702	.8810 / 41.4081	1.4294	.8609 / 42.5002	1.2176	.6808 / 42.4959
DN	(1*)	1.2702	–	–	–	–	–
CBC	(2)	1.8120	.9730 / 76.4071	2.9126	.8884 / 92.7365	2.4538	.7044 / 92.7399
PF	(2)	1.4139	1.0091 / 41.1094	2.0604	1.2238 / 64.1713	2.4531	1.5846 / 61.3726
CSF- $B_{wp}$	(3)	1.3031	.8757 / 43.8770	1.4833	.8122 / 44.6721	1.2641	.6472 / 44.7587
EPF	(3)	1.3705	– / 21	–	–	–	–
ProjSFM-BA	(3*)	–	–	–	–	–	.64 / 41.5

<sup>†</sup>Models: (1) general rank-4 factors (no mean column); (2) affine SFM; (3) Euclidean (weak-perspective) SFM; (3\*) projective SFM with Euclidean bundle adjustment; (1\*) penalty terms imposed on model (1) to favor orthogonality of camera axes. Error values are given in pixels.

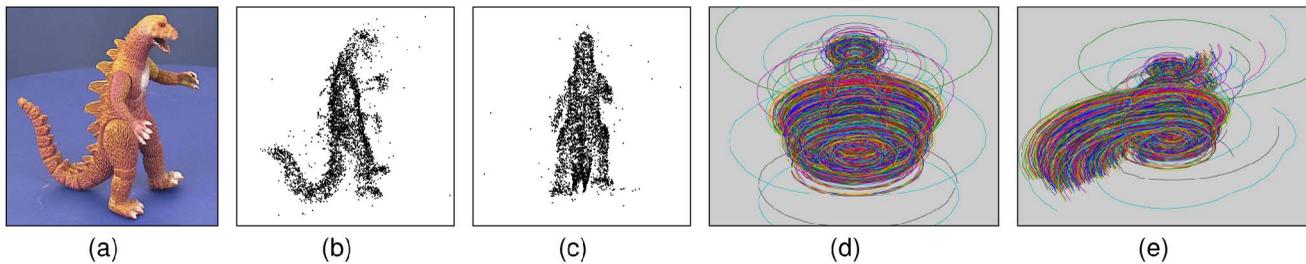


Fig. 4. Results on the 4,983 point set (90.8 percent missing) of the dinosaur data set: (a) one of the 36 images of the sequence, (b)-(c) side and frontal views of the euclidean 3D shape recovered by CSF- $B_{wp}$ , and (d) 2D point tracks reconstructed by CSF- $B_{wp}$  and (e) by the best affine methods.

500 for the other methods). CSF- $B_{af}$  and CSF- $B_{wp}$  are also very fast. Typically, CSF- $B_{wp}$  performs more iterations than CSF- $B_{af}$ —for better convergence on long sequences (large  $T$ ), we run CSF- $B_{wp}$  in a coarse-to-fine manner, increasing  $d \in \{0.1T, 0.2T, 0.3T\}$ . We also note that the Wiberg method does not scale well for application on large matrices due to the size of the system of equations it has to solve in each iteration.

Our experiments also considered the complete dinosaur data set with 4,983 tracks (of which 2,300 are defined on only two images).<sup>1</sup> The dinosaur sequence is arguably the most popular data set used to evaluate rigid SFM algorithms, with results published on different subsets of its 2D point tracks. For comparison against previously published results, we consider a subset of 2,683 tracks (points tracked in at least three images) and the subset of 319 tracks described above.

Results of the methods above are also compared to those of the following algorithms: Damped-Newton (DN) [6], the Deviation Parameter (DP) for subspace constraints [19], Minimal Missing Elements (MME) [7], Camera Basis (subspace) Constraints (CBC) [28], and the Euclidean Power-Factorization (EPF) method in [21]. Finally, we also compare the result of CSF- $B_{wp}$  to that obtained with projective SFM followed by euclidean bundle adjustment (ProjSFM-BA) [22]. Because the ground-truth 3D shape and cameras are not available, we compare RMSE values for the

reconstructed  $\mathbf{W}$  and also the mean/maximum 2D reprojection errors (in pixels) of the best solutions provided by each method. For this numerical comparison on the relatively short dinosaur sequence ( $T = 36$ ), our CSF methods were both initialized (deterministically) with full DCT bases.

In Table 4, each incomplete row contains only the results reported in the original publication due to the algorithm not being available in our experiments. Because these algorithms use factorization models with different numbers of degrees of freedom, we define three main groups of methods. In this case, models with more degrees of freedom are expected to be associated with lower RMSE and 2D error values. However, we note that CSF- $B_{wp}$  yielded lower 2D reprojection errors compared to the affine SFM methods. In the dinosaur data set, the performance of CSF- $B_{wp}$  is very close to that of the projective SFM method with bundle adjustment, despite the nonnegligible perspective distortion in the data. Figs. 4b and 4c show the 3D shape recovered by CSF- $B_{wp}$ .

Information on the ground-truth motion is actually available because we know the dinosaur is on a turntable. Thus, the reconstructed  $\mathbf{W}$  must describe 2D point trajectories that are all concentric ellipses. Indeed, the result of CSF- $B_{wp}$  indicates correct motion recovery, Fig. 4d. The best affine methods in our comparison (CSF- $B_{af}$ , LM-S, LM- $S_{GN}$ , and Wiberg) overfit noise and outliers and recover an incorrect motion pattern, Fig. 4e. This problem is inherent

1. Available at <http://www.robots.ox.ac.uk/~vgg/>.

TABLE 5  
Results of Rigid SFM Methods on the Teddy Bear Data Set

Method	(model)	RMSE	mean/max 2D error
CSF- $B_{af}$	(1)	.4925	.4527 / 7.4998
LM-S, LM- $S_{GN}$	(1)	.4925	.4527 / 7.4998
DP	(1)	.5911	.4977 / 13.8324
CSF- $B_{af}$	(2)	.6174	.5689 / 10.7385
Wiberg	(2)	.6174	.5689 / 10.7385
CBC	(2)	.7487	.6478 / 14.1339
PF	(2)	1.2125	.8649 / 27.5116
CSF- $B_{wp}$	(3)	.6310	.5792 / 10.6524

†Models are as in Table 4. Error values are given in pixels.

in all affine SFM methods (incorrect 2D trajectories are also shown in [6] and [8]). Also note this fact is independent of the euclidean upgrade step, which does not affect reprojection.

The robust CBC method uses random sampling to tolerate outliers when computing subspace constraints and triangulating 3D points. Thus, CBC yields higher RMSEs as compared to the other methods. Table 4 shows that 2D reprojection error decreases as more point tracks are available for outlier/inlier identification—we computed camera constraints from triplets of frames presenting at least 15 tracks in common. In [28], CBC was found to provide the best solutions among the main subspace constraint methods. Here, error values are higher than those in [28] because we did not remove outliers before the comparison against nonrobust methods. We note that our CSF methods do not explicitly address the problem of outliers and may fail if the input data present gross measurement errors. Future work will investigate the use of robust error terms [9], [12] in our cost functions.

In another experiment, we applied the Kanade-Lucas-Tomasi (KLT) feature tracker [27] on a 200-image *teddy bear* sequence [28]. We selected 806 point tracks appearing in at least 10 images. Due to occlusion, the resulting matrix  $\mathbf{W} \in \mathbb{R}^{400 \times 806}$  is missing 88.6 percent of its entries. For numerical comparison, both CSF methods considered a full DCT basis. CSF- $B_{wp}$  adopted a coarse-to-fine strategy with  $d \in \{0.1T, 0.2T, 0.3T, \dots, T\}$  increased each time convergence to a coarser solution was detected. Table 5 gives the 2D errors obtained with the algorithms described above on our teddy bear data set. CSF- $B_{af}$  provides equivalent or better solutions as compared to the other methods. Also, CSF- $B_{wp}$  directly reconstructs the euclidean 3D shape with only a slight increase in 2D error—its camera model captures less of the noise in the data. Fig. 5 shows the 3D shape reconstructed with CSF- $B_{wp}$ . On this long sequence ( $T = 200$ ) with smooth motion, the difference in the results of CSF- $B_{wp}$  with a full DCT basis ( $d = T$ ) and with  $d = 0.3T$  is very small,  $e_R = 0.0175$  and  $e_S = 0.0229$  (here,  $e_S$  is normalized by  $\sigma(\mathbf{S})$  as in (21)). With  $d = 0.3T$ , the runtime of CSF- $B_{wp}$  drops from 27.3 to 3.6 minutes (4 to 0.7 minutes for CSF- $B_{af}$ ). Wiberg took, on average, 4.2 hours.

### 5.3 Computing Nonrigid SFM

First, we evaluate our nonrigid SFM algorithm on complete data sets with known 3D shapes for each frame, also

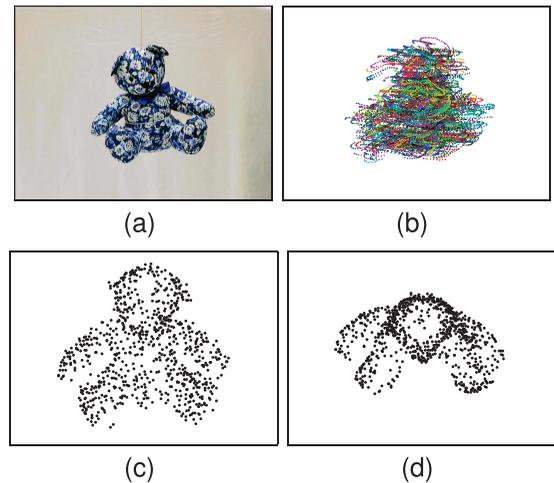


Fig. 5. Result of CSF- $B_{wp}$  on 806 point tracks (88.6 percent missing) of the teddy bear sequence: (a) image in the sequence, (b) input 2D tracks, (c)-(d) frontal and top views of the recovered euclidean 3D shape.

simulating missing data and noise. Then, we present results on real data sets with and without occlusion. The number of frames ( $T$ ) and the number of point tracks ( $n$ ) are indicated as ( $T/n$ ) after a data set's name.

We start with the motion capture sequences: *drink* (1,102/41), *pickup* (357/41), *yoga* (307/41), *stretch* (370/41), and *dance* (264/75) used in [2]; *face1* (74/37) of [24]; *face2* (316/40) and *walking* (260/55) of [30]. We also use the synthetic bending *shark* (240/91) of [30]. Note that a different *shark* data set appears in [2]. We use only the original one in [30].

To allow for comparison against the results reported in [2], we used the same procedure and error metrics therein. For each data set, the complete 2D point trajectories in  $\mathbf{W}$  are obtained by applying an orthographic projection on the sequence of 3D shapes. Because the solution of nonrigid SFM methods is defined up to an arbitrary  $3 \times 3$  rotation, we compute a single rotation that best aligns all reconstructed and original 3D shapes. Let  $e_{tj}$  be the reconstruction error (i.e., euclidean distance) for the  $j$ th 3D point of frame  $t$ . We then compute a normalized mean 3D error over all points and frames,

$$e_{3D} = \frac{1}{\sigma_e T n} \sum_{t=1}^T \sum_{j=1}^n e_{tj}, \quad \sigma_e = \frac{1}{T} \sum_{t=1}^T \sigma(\mathbf{S}_t),$$

with  $\sigma(\cdot)$  as in (21) and  $\mathbf{S}_t \in \mathbb{R}^{3 \times n}$  the original 3D shape in frame  $t$ . The first four motion capture sequences have artificial rotations applied on them. We thus compare original and estimated rotations using  $e_R$  as in (22).

Table 6 compares the performance of our 3D shape trajectory approach, CSF- $B_{nr}$ , against four state-of-the-art, nonrigid SFM methods:

1. the shape basis constraints (XCK) method<sup>2</sup> [32],
2. the algorithm modeling 3D shape using probabilistic principal component analysis (EM-PPCA) [30],

<sup>2</sup>We did not implement the XCK algorithm. Its results on the first five sequences are reproduced from [2]. Results on the shark, face, and walking data sets are given in [30], with a different error metric, and found to be significantly inferior to those of EM-PPCA.

TABLE 6  
Performance of Nonrigid SFM Methods on Synthetic and Motion Capture Data

Dataset	XCK		EM-PPCA		MP		PTA		CSF-B <sub>nr</sub>		
	$e_R$	$e_{3D}$	$e_R$	$e_{3D}$	$e_R$	$e_{3D}$	$e_R$	$e_{3D}$ (K)	$e_R$	$e_{3D}$ (K)	initial $e_{3D}$ (K)
Drink	.3359	3.5186	.2906	.3393	.2859	.4604	<b>.0058</b>	<b>.0250</b> (13)	<b>.0055</b>	<b>.0223</b> (6)	.0854 (6)
Pick-up	.4687	3.3721	.4277	.5822	.2506	.4332	<b>.1549</b>	<b>.2369</b> (12)	<b>.1546</b>	<b>.2301</b> (6)	.2685 (6)
Yoga	1.2014	7.4935	.8089	.8097	.8711	.8039	<b>.1059</b>	.1625 (11)	<b>.1021</b>	<b>.1467</b> (7)	.1528 (7)
Stretch	.9489	4.2415	.7594	1.1111	.8174	.8549	.0549	.1088 (12)	<b>.0489</b>	<b>.0710</b> (8)	.0966 (8)
Dance	–	2.9962	–	.9839	–	<b>.2639</b>	–	.2958 (5)	–	<b>.2705</b> (2)	.3259 (2)
Face1	–	–	–	<b>.0434</b>	–	.0734	–	.1247 (3)	–	.0637 (5)	.1487 (3)
Face2	–	–	–	<b>.0329</b>	–	<b>.0357</b>	–	.0444 (5)	–	<b>.0363</b> (3)	.0451 (3)
Shark	–	–	–	.0501	–	.1571	–	.1796 (9)	–	<b>.0081</b> (3)	.3195 (3)
Walking	–	–	–	.4917	–	.5607	–	.3954 (2)	–	<b>.1863</b> (2)	.6823 (2)

For the related PTA and CSF-B<sub>nr</sub> methods, factorization rank is also indicated by the value of  $K$  in parenthesis.

3. the Metric Projections (MP) method<sup>3</sup> [24], and
4. the DCT-based 3D point trajectory approach [2].

Following the methodology in [2], we ran the algorithms with different values of  $K \in \{2, 3, \dots, 13\}$ , reporting the best result. Table 6 also shows the value of  $K$  for the best solutions obtained with PTA and CSF-B<sub>nr</sub> for comparison. We also report the initial error of CSF-B<sub>nr</sub> because it shows the error that PTA would provide with the same  $K$  and rotations in  $\mathbf{D}$ . In all runs, CSF-B<sub>nr</sub> had the number of DCT basis set to  $d = 0.1T$ , except for the two face data sets on which we set  $d = \frac{T}{3}$  due to the presence of higher frequency deformations.

Our results as shown in Table 6 are consistently similar to or better than the best results provided by the other methods on each data set. As compared to PTA, CSF-B<sub>nr</sub> computes better solutions and at a lower rank  $r = 3K + 1$  by more efficiently using higher DCT frequency components. Furthermore, our simple strategy of iterating over  $K$  while computing factor  $\mathbf{D}$  also provides better euclidean camera estimates than plain PTA.

EM-PPCA, MP, and CSF-B<sub>nr</sub> provide comparable results on the face2 data set, which has mostly rigid motion with high-frequency deformation seen on the lower lips and chin. In this case, CSF-B<sub>nr</sub> requires at least 75 percent of all DCT components ( $d = 0.75T$ ) to provide a mean 3D error (0.0328) that is smaller than that of EM-PPCA. PTA is not capable of modeling this high-frequency deformation and recovers a mostly rigid mouth, Fig. 6. On the much shorter face1 sequence, CSF-B<sub>nr</sub> obtains an error improvement of only 0.0012 with a larger  $d$ . Setting  $d > 0.1T$  resulted in no significant improvement to the CSF-B<sub>nr</sub> solutions on the other data sets.

CSF-B<sub>nr</sub> is the only method that can accurately reconstruct the deformation of the bending shark. Furthermore, with a full DCT basis, the resulting  $e_{3D}$  is negligible (0.00004) and perfect reconstruction is achieved. Fig. 7 shows the three best shark reconstructions of Table 6. XCK

3. Two MP methods are presented in [24]. Here we experiment with the method using a generic model of deformable shapes. The supplementary file, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>, shows comparative results against MP with a specialized articulated model that requires the columns of  $\mathbf{W}$  to be initially grouped as corresponding to separate object parts.

fails on this data set due to the bending shark presenting a 2D (degenerate) deformation mode [30].

The motion capture sequences containing highly articulated bodies also show the superiority of CSF-B<sub>nr</sub> and PTA compared to XCK, EM-PPCA, and MP. On the smoothly deforming shapes in the drink, pickup, and yoga sequences (used in [2]), the improvement offered by CSF-B<sub>nr</sub> over PTA is marginal. On the other hand, the results on the more difficult walking sequence of [30] (Fig. 8) highlight the advantages of using our 3D shape trajectory approach (no artificial rotation was added to this sequence). MP provides comparable results with CSF-B<sub>nr</sub> on the dance sequence. The supplementary file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>, includes additional examples.

To simulate missing data in the shark and walking data sets, we randomly discard  $\rho$  percent of the 2D entries in  $\mathbf{W}$ . Before applying CSF-B<sub>nr</sub> (with  $K$  as in Table 6), we compute  $\mathbf{D}$  and  $\mathbf{t}$  by first using CSF-B<sub>af</sub> to reconstruct the complete 2D point trajectories in  $\mathbf{W}$ . CSF-B<sub>af</sub> was run with  $d = 0.25T$  and rank  $r = 7$ . Letting  $\mathbf{W}_0$  be the complete matrix, we normalize the 2D reconstruction error for the incomplete  $\mathbf{W}$  by  $\sigma(\mathbf{W}_0)$  as in (21).

On the smooth shark deformation, results are visually similar to those in Fig. 7, with  $\rho$  up to 95 percent (see images in the supplementary file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>). Of 10 runs with  $\rho = 95\%$ , the average (maximum) 2D reconstruction error for  $\mathbf{W}$  was 0.0015 (0.0029). The average (maximum) 3D error after running CSF-B<sub>nr</sub> was 0.0163 (0.0488). On the walking

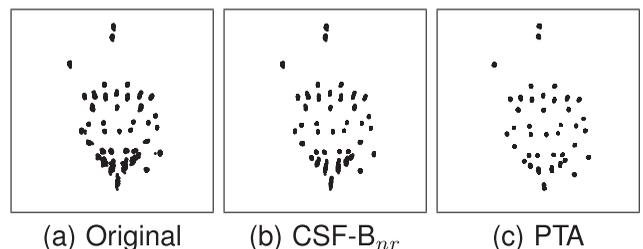


Fig. 6. Overlay of 316 unrotated 3D shapes in the face2 sequence. Deformation is seen predominantly on the lower-lips and chin. CSF-B<sub>nr</sub> can capture this high-frequency deformation better than PTA.

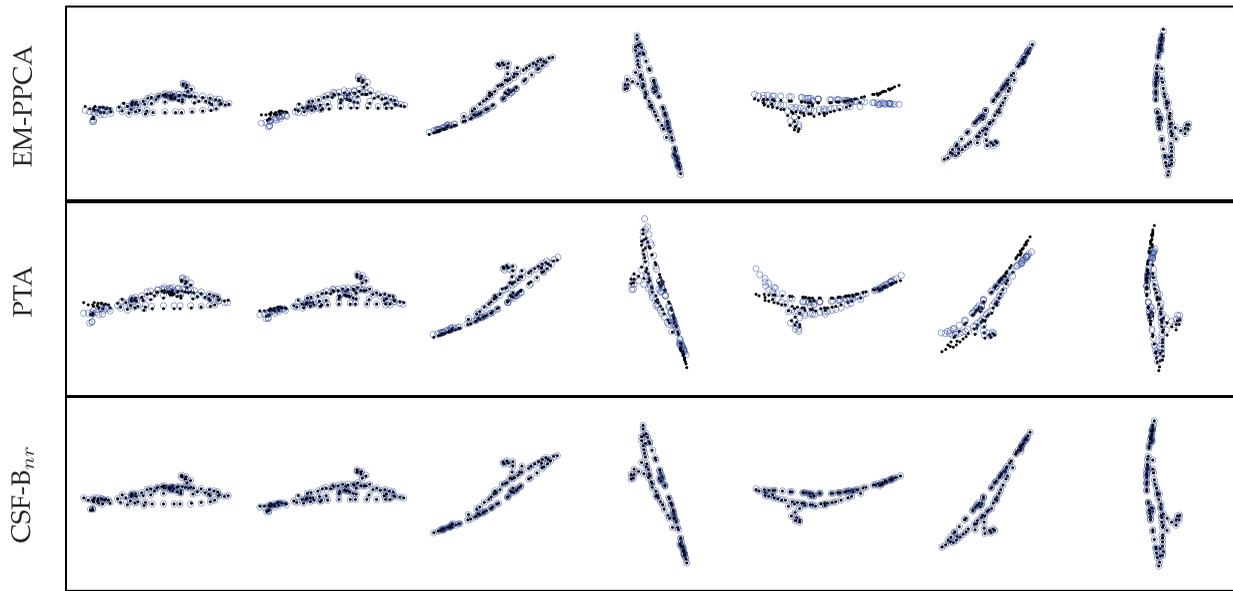


Fig. 7. Results on the bending shark sequence. Reconstructed 3D shapes (blue circles) are shown against the original 3D data (dark dots). Frames 20, 50, 80, 115, 148, 175, and 200 are displayed above.

sequence, results with  $\rho = 75\%$  are still visually similar to those in Fig. 8. After 10 runs, the average (maximum) 2D and 3D errors were 0.0508 (0.0548) and 0.2063 (0.3910), respectively. These results on incomplete data are still better than those of XCK, EM-PPCA, MP, and PTA on the complete shark and walking data sets.

On the complete face2 data set, the performances of EM-PPCA, MP, PTA, and CSF- $B_{nr}$  are more similar than on the other sequences. We thus analyzed their average  $e_{3D}$  with different levels of random occlusion and noise (simulated independently) on  $\mathbf{W}$  of face2. EM-PPCA and MP handle cases of missing data by adopting an alternation approach as in PowerFactorization. PTA does not handle occlusions and was tested with added noise only. All methods were

run with their best parameter  $K$  for the complete data ( $K = 5$  for EM-PPCA, MP, PTA;  $K = 3$  for CSF- $B_{nr}$ ). After random occlusion, we ensured that  $\mathbf{W}$  had at least  $3K_{max} + 1$  entries ( $K_{max} = 5$ ) in each row and column. On the high-frequency deformation and motion of face2 (combined in the columns of  $\mathbf{W}$ ), the initial CSF- $B_{af}$  step of CSF- $B_{nr}$  reconstructed the 2D point trajectories in  $\mathbf{W}$  with  $d = 0.5T$  and  $r = 5$ . Results were averaged over 100 trials and are shown in Fig. 9.

In our trials on face2 with missing data, EM-PPCA was unstable and presented both small and very large errors for random occlusions  $\rho$  as low as 10 percent. Thus, Fig. 9a shows the results of CSF- $B_{nr}$  in comparison to those of MP, which was shown in [24] to outperform EM-PPCA in a

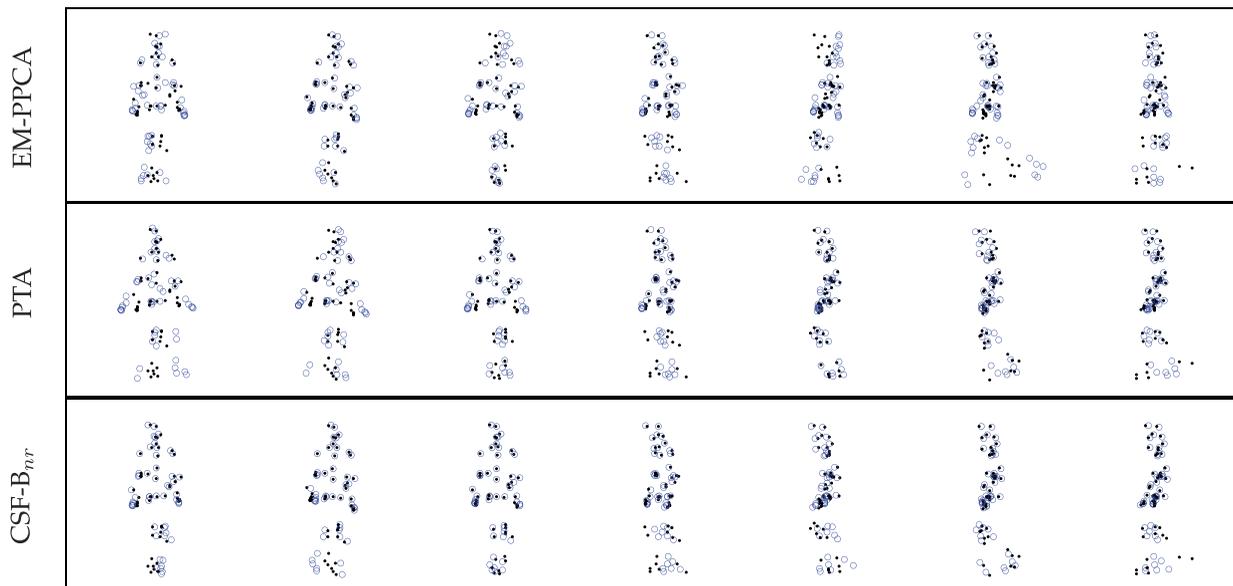


Fig. 8. Results on the walking sequence. Reconstructed 3D shapes (blue circles) are shown against the original 3D data (dark dots). Frames 34, 74, 122, 160, 198, 223, and 255 are displayed above.

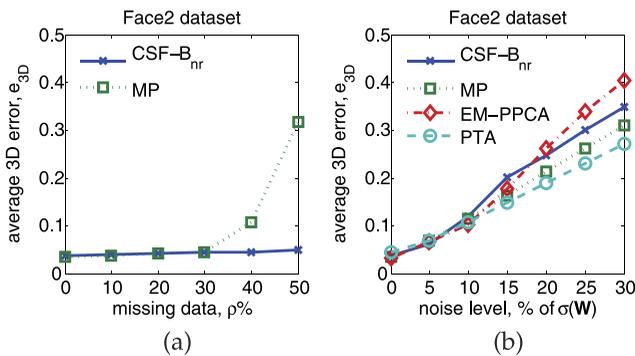


Fig. 9. Reconstruction errors on the face2 sequence with missing data and noise simulated independently (averages over 100 trials).

similar experiment on face2. Note that the average  $e_{3D}$  of MP begins to increase with  $\rho$  above 30 percent, while that of CSF- $B_{nr}$  presents almost no variation over all the tested levels of random occlusion.

For levels of added Gaussian noise below  $0.15\sigma(\mathbf{W})$  (15 percent), the increase in  $e_{3D}$  is similar for all four methods, Fig. 9b. At 15 percent and above, small differences in performance are observed but the reconstructions still degrade only gradually with the level of noise. Because PTA largely oversmooths the deformation of face2, its performance suffers a smaller penalty as compared to that of the other methods. At high-noise levels, correct 3D shape reconstruction seems to require additional information on the nature of the 3D shapes and noise [12].

The average runtimes of the algorithms on the face2 trials at 5 percent noise level were, in minutes: 2.52 (EM-PPCA), 0.34 (MP), 0.05 (PTA), and 0.40 (CSF- $B_{nr}$ , of which 0.15 was spent to compute  $\mathbf{D}$ ). These times were similar to those observed for the algorithms on the original data. PTA is fast due to its use of SVD. Runtimes in minutes on face2 trials with 30 percent missing data were: 2.85 (EM-PPCA), 3.66 (MP), and 2.29 (CSF- $B_{nr}$ ). This runtime of CSF- $B_{nr}$  includes time spent running CSF- $B_{af}$  (0.82) and computing  $\mathbf{D}$  (0.19). At above 20 percent random occlusion, CSF- $B_{nr}$  was faster than MP and EM-PPCA.

We also applied CSF- $B_{nr}$  to the (complete) real data set *cubes* (200/14) of [2]. With  $K = 2$  and  $d = 0.1T$ , the solution of CSF- $B_{nr}$  has mean (maximum) 2D reprojection error of 0.4958 (2.0672) pixel. The solution of PTA has an error of 1.6589 (4.8602) also with  $K = 2$ . Images and additional results are given in the supplementary file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>.

Finally, an application of our nonrigid SFM method is in the interpretation of the facial expression component of sign languages from video [10]. In this case, head rotation and hand gesticulation often cause the occlusion of facial features, leading to incomplete 2D point tracks. We now consider a 115-image (4 seconds long) face close-up sequence of an American Sign Language (ASL) sentence. Facial landmarks were manually annotated in each image when visible. The resulting matrix  $\mathbf{W} \in \mathbb{R}^{230 \times 77}$  is missing 17.4 percent of its data and has small magnitude noise due to annotation errors caused by partial occlusion of facial features and motion blur in the video images. Fig. 10 shows six example images and their respective 3D face shapes recovered using CSF- $B_{nr}$  ( $K = 2$  and  $d = 0.5T$ , with the initial CSF- $B_{af}$  step run as for face2 above). Figures with the results of EM-PPCA and MP are given in the supplementary file, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2011.50>. For all tested values of  $K \in \{2, 3, \dots, 13\}$ , the result obtained with EM-PPCA presented visibly large 3D reconstruction errors for a number of facial points that are occluded frequently in each sequence (e.g., the lateral contours of the nose and the face). The results of CSF- $B_{nr}$  and MP are visually similar and indicate correct recovery of pose and nonrigid 3D shape despite the occurrence of occlusion.

## 6 CONCLUSION

In this paper, we have addressed the classical computer vision problems of rigid and nonrigid SFM with occlusion. We started by assuming that the columns of the input data matrix  $\mathbf{W}$  describe smooth 2D point trajectories over time.

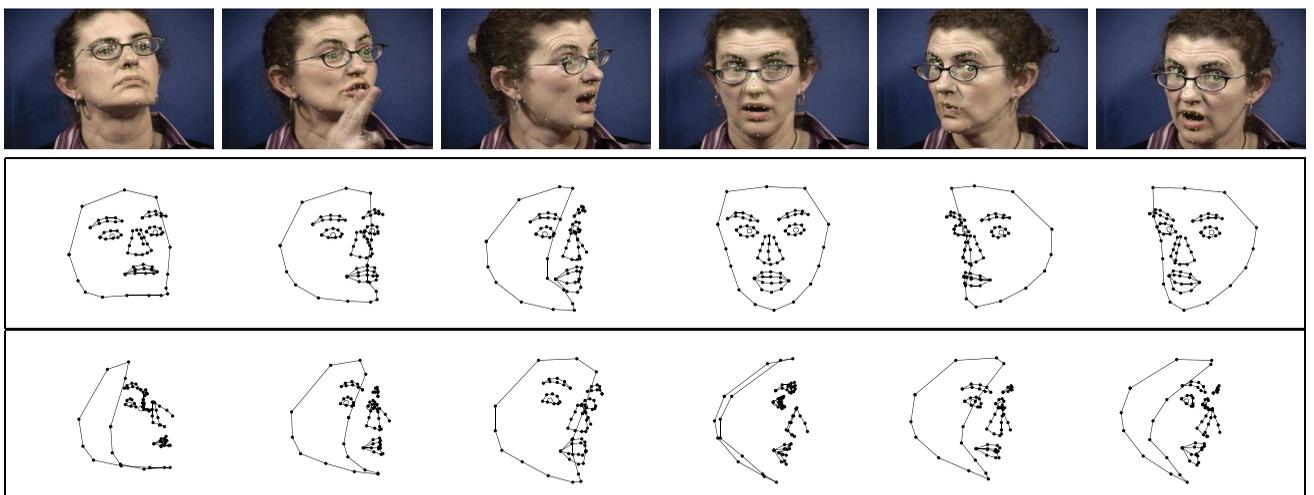


Fig. 10. Results of CSF- $B_{nr}$  on the ASL sequence (77 points, 17.4 percent missing data): (top) six out of 115 images with annotated facial landmarks in green; (middle, bottom) two orthogonal views of the recovered 3D shapes.

This assumption is equivalent to considering that 2D observations in  $\mathbf{W}$  are obtained by a single camera moving smoothly around a rigid structure. In nonrigid SFM, our assumption also requires the structure to deform only smoothly over time.

Our main contributions are two-fold: 1) We provide new models for the smooth time-trajectories of camera and deformable shape with a compact parameterization in the DCT domain, and 2) we derive a family of efficient Column Space Fitting methods to estimate such trajectories while tolerating cases in which  $\mathbf{W}$  is presented with missing data.

In rigid SFM, we consider a weak-perspective camera model from the outset and directly reconstruct euclidean 3D shape without requiring postprocessing steps. Our results on synthetic and real SFM data sets with noise and high percentages of missing data were positively compared to the state of the art.

In nonrigid SFM, we propose a novel 3D shape trajectory approach that solves for the deformable structure as the trajectory of a single point in an implicitly defined linear shape space. A comparison against state-of-the-art algorithms shows that our method can better model complex articulated deformation with higher frequency DCT components while still maintaining the low-rank factorization constraint. We also demonstrate that our nonrigid SFM algorithm can tolerate high percentages of missing data in the input matrix  $\mathbf{W}$  with only a small penalty in 3D reconstruction accuracy.

Future work will investigate the integration of our weak-perspective camera model into our nonrigid SFM approach. We will also consider the automatic selection of the number of elements in the DCT and shape bases using regularization terms that balance the trade-off between higher model complexity ( $r$ ,  $K$ , and  $d$ ) and smaller fitting error [23], [25], [30]. Our CSF algorithms may also benefit from a robust error term in order to better tolerate the presence of outliers in  $\mathbf{W}$ . Finally, nonlinear models for the column space of  $\mathbf{W}$  shall be considered.

## ACKNOWLEDGMENTS

The authors thank the reviewers for their constructive comments. This research was supported in part by the US National Science Foundation (NSF), grant 0713055, and the US National Institutes of Health, grant R01 DC 005241.

## REFERENCES

- [1] I. Akhter, Y. Sheikh, and S. Khan, "In Defense of Orthonormality Constraints for Nonrigid Structure from Motion," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 1534-1541, 2009.
- [2] I. Akhter, Y.A. Sheikh, S. Khan, and T. Kanade, "Nonrigid Structure from Motion in Trajectory Space," *Proc. Neural Information Processing Systems*, Dec. 2008.
- [3] A. Bartoli, V. Gay-Bellile, U. Castellani, J. Peyras, S. Olsen, and P. Sayd, "Coarse-to-Fine Low-Rank Structure-from-Motion," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, no. 1, pp. 1-8, 2008.
- [4] D. Bertsekas, *Nonlinear Programming*. Athena Scientific Press, 1999.
- [5] C. Bregler, A. Hertzmann, and H. Biermann, "Recovering Non-Rigid 3d Shape from Image Streams," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 2, pp. 690-696, 2000.
- [6] A.M. Buchanan and A.W. Fitzgibbon, "Damped Newton Algorithms for Matrix Factorization with Missing Data," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 2, pp. 316-322, 2005.
- [7] P. Chen and D. Suter, "Recovering the Missing Components in a Large Noisy Low-Rank Matrix: Application to SFM," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 26, no. 8, pp. 1051-1063, Aug. 2004.
- [8] P. Chen, "Optimization Algorithms on Subspaces: Revisiting Missing Data Problem in Low-Rank Matrix," *Int'l J. Computer Vision*, vol. 80, no. 1, pp. 125-142, 2008.
- [9] F. de la Torre and M. Black, "Robust Principal Component Analysis for Computer Vision," *Proc. Int'l Conf. Computer Vision*, vol. 1, pp. 362-369, 2001.
- [10] L. Ding and A. Martinez, "Modelling and Recognition of the Linguistic Components in American Sign Language," *Image and Vision Computing*, vol. 27, no. 12, pp. 1826-1844, 2009.
- [11] D. Forsyth and J. Ponce, *Computer Vision: A Modern Approach*. Prentice Hall, 2003.
- [12] J. Fortuna and A.M. Martinez, "Rigid Structure from Motion from a Blind Source Separation Perspective," *Int'l J. Computer Vision*, 2010.
- [13] G.H. Golub and C.F. Van Loan, *Matrix Computations*. Johns Hopkins Studies in Math. Sciences, The Johns Hopkins Univ. Press, 1996.
- [14] N. Guilbert, A. Bartoli, and A. Heyden, "Affine Approximation for Direct Batch Recovery of Euclidian Structure and Motion from Sparse Data," *Int'l J. Computer Vision*, vol. 69, no. 3, pp. 317-333, 2006.
- [15] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*. Cambridge Univ. Press, 2003.
- [16] R. Hartley and F. Schaffalitzky, "Powerfactorization: 3D Reconstruction with Missing or Uncertain Data," *Proc. Australia-Japan Advanced Workshop Computer Vision*, 2003.
- [17] D. Jacobs, "Linear Fitting with Missing Data for Structure-from-Motion," *Computer Vision and Image Understanding*, vol. 82, no. 1, pp. 57-81, 2001.
- [18] A. Jain, "A Sinusoidal Family of Unitary Transforms," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 1, no. 4, pp. 356-365, Oct. 1979.
- [19] H. Jia and A.M. Martinez, "Low-Rank Matrix Fitting Based on Subspace Perturbation Analysis with Applications to Structure from Motion," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 31, no. 5, pp. 841-854, May 2009.
- [20] J. Magnus and H. Neudecker, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, second ed. Wiley, 1999.
- [21] M. Marques and J. Costeira, "Estimating 3D Shape from Degenerate Sequences with Missing Data," *Computer Vision and Image Understanding*, vol. 113, no. 2, pp. 261-272, 2009.
- [22] D. Martinec and T. Pajdla, "Structure from Many Perspective Images with Occlusions," *Proc. European Conf. Computer Vision*, pp. 355-369, 2002.
- [23] S.I. Olsen and A. Bartoli, "Implicit Non-Rigid Structure-from-Motion with Priors," *J. Math. Imaging Vision*, vol. 31, nos. 2/3, pp. 233-244, 2008.
- [24] M. Paladini, A. Del Bue, M. Stošić, M. Dodig, J. Xavier, and L. Agapito, "Factorization for Non-Rigid and Articulated Structure Using Metric Projections," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 2898-2905, 2009.
- [25] M. Pollefeys, F. Verbiest, and L. Van Gool, "Surviving Dominant Planes in Uncalibrated Structure and Motion Recovery," *Proc. European Conf. Computer Vision*, vol. 2, pp. 837-851, 2002.
- [26] V. Rabaud and S. Belongie, "Rethinking Nonrigid Structure from Motion," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 1-8, 2008.
- [27] J. Shi and C. Tomasi, "Good Features to Track," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, no. 1, pp. 593-600, 1994.
- [28] J.-P. Tardif, A. Bartoli, M. Trudeau, N. Guilbert, and S. Roy, "Algorithms for Batch Matrix Factorization with Application to Structure-from-Motion," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, no. 1, pp. 1-8, 2007.
- [29] C. Tomasi and T. Kanade, "Shape and Motion from Image Streams under Orthography: A Factorization Method," *Int'l J. Computer Vision*, vol. 9, no. 2, pp. 137-154, Nov. 1992.
- [30] L. Torresani, A. Hertzmann, and C. Bregler, "Nonrigid Structure-from-Motion: Estimating Shape and Motion with Hierarchical Priors," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, no. 5, pp. 878-892, May 2008.
- [31] T. Wiberg, "Computation of Principal Components When Data Is Missing," *Proc. Second Symp. Computational Statistics*, pp. 229-236, 1976.

- [32] J. Xiao, J. Chai, and T. Kanade, "A Closed Form Solution to Non-Rigid Shape and Motion Recovery," *Int'l J. Computer Vision*, vol. 67, no. 2, pp. 233-246, 2006.
- [33] J. Xiao and T. Kanade, "Non-Rigid Shape and Motion Recovery: Degenerate Deformations," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 668-675, 2004.
- [34] J. Yan and M. Pollefeys, "A Factorization-Based Approach for Articulated Nonrigid Shape, Motion and Kinematic Chain Recovery from Video," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, no. 5, pp. 865-877, May 2008.



**Paulo F.U. Gotardo** received the BSc and MSc degrees in Informatics from Universidade Federal do Paraná (UFPR), Brazil, in 2000 and 2002, respectively, and the PhD degree in electrical and computer engineering from The Ohio State University (OSU) in 2010. Currently, he is a postdoctoral researcher with the Computational Biology and Cognitive Science Lab (CBCSL) in OSU's Department of Electrical and Computer Engineering. Prior to joining CBCSL,

he was with the IMAGO Research Group at UFPR and with the Signal Analysis and Machine Perception Laboratory at OSU. He also worked on human-computer interaction (HCI) in immersive environments at OSU's Advanced Computing Center for the Arts and Design. His research interests include machine vision, learning, and HCI. He is a member of the IEEE.



**Aleix M. Martinez** is an associate professor in the Department of Electrical and Computer Engineering at The Ohio State University (OSU), where he is the founder and director of the Computational Biology and Cognitive Science Lab. He is also affiliated with the Department of Biomedical Engineering and with the Center for Cognitive Science. Prior to joining OSU, he was affiliated with the Electrical and Computer Engineering Department at Purdue

University and with the Sony Computer Science Lab. He is an associate editor of the *IEEE Transactions on Pattern Analysis and Machine Intelligence* and of *Image and Vision Computing*, and has served as an area chair for CVPR, ICPR, and F&G. His areas of interest are learning, vision, linguistics, and their interaction. He is a senior member of the IEEE.

▷ **For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).**