

# Estimating the Pose of a 3D Sensor in a Non-Rigid Environment

Adrien Bartoli

CNRS - LASMEA, France — [Adrien.Bartoli@gmail.com](mailto:Adrien.Bartoli@gmail.com)

## Abstract

*Estimating the pose of an imaging sensor is a central research problem. Many solutions have been proposed for the case of a rigid environment. In contrast, we tackle the case of a non-rigid environment observed by a 3D sensor, which has been neglected in the literature. We represent the environment as sets of time-varying 3D points explained by a low-rank shape model, that we derive in its implicit and explicit forms. The parameters of this model are learnt from data gathered by the 3D sensor. We propose a learning algorithm based on minimal 3D non-rigid tensors that we introduce. This is followed by a Maximum Likelihood non-linear refinement performed in a bundle adjustment manner. Given the learnt environment model, we compute the pose of the 3D sensor, as well as the deformations of the environment, that is, the non-rigid counterpart of pose, from new sets of 3D points. We validate our environment learning and pose estimation modules on simulated and real data.*

## 1. Introduction

Aligning 3D views – sets of 3D points – gathered by a 3D sensor, such as a calibrated stereo rig, is important for constructing comprehensive 3D models of the environment or updating the position of a mobile imaging sensor. When the environment is rigid, the 3D views are related by rigid Euclidean transformations. Many approaches have been proposed to compute these transformations, *e.g.* [8]. Aligning 3D views is one of the building blocks of hierarchical approaches to Structure-From-Motion. However, the assumption of rigidity is violated in many cases of interest, for instance a garment deforming as a person moves. The alignment problem is then particularly challenging because a different shape is observed in each 3D view.

A large body of work has been done in the medical imaging community, see [7] for a review, but with the aim of estimating dense deformation fields from dense, often voxel-based, reconstructions. Dealing with non-rigid scenes coming from single-camera footage has received an increasing attention over the last few years. The problem is highly challenging since both the cameras and the non-rigid shape have to be recovered. A major step forwards for such cases was made by Bregler *et al.* [4, 10] and Brand [3]. Build-

ing on the work of [2, 9], they developed and demonstrated factorization of images of non-rigid scenes, where the non-rigidity was represented as a linear combination of basis shapes. It is shown in [5] how the constraints coming from two synchronized cameras can be incorporated into non-rigid factorization.

We tackle the problem of computing the pose of a 3D sensor with respect to a non-rigid scene, that we represent using the low-rank shape model used in non-rigid factorization methods. Most previous work, *e.g.* [2, 3, 4, 5, 10, 13] use the weak perspective camera model. In contrast, we do not specify a camera model, since we directly consider 3D views. We assume that spatial and temporal point correspondences are established. Pose estimation in a non-rigid environment raises two main problems. First, one has to define the meaning of non-rigid pose. This problem is tackled in [14] in a different context, where the notions of ‘shape average’ and ‘moving average’ are introduced. One benefit of using the low-rank shape model is that the ‘true’ camera pose is recovered. Second, contrarily to classical model-based pose estimation in a rigid environment, a prior model of the non-rigid environment is not available in many cases. We propose to learn this model from a collection of unregistered 3D views gathered by the 3D sensor. Once this learning stage has been passed, our non-rigid pose estimator can be launched.

We bring the following contributions. First, §3, we state the implicit and explicit low-rank shape models, and state the notion of pose in this context. Second, §4, we propose algorithms to learn the non-rigid environment. The implicit model parameters are learnt using a factorization technique, while for the explicit model, we use what we call *minimal 3D non-rigid tensors*. Third, §5, we show how the pose of the 3D sensor can be computed with respect to the learnt model while the environment is moving and deforming. Experimental results on simulated and real data are reported in §6. We give our conclusions in §7.

## 2. Notation

Matrices are written in sans-serif fonts, *e.g.*  $R$ , and vectors using bold fonts, *e.g.*  $\mathbf{x}$ . The  $n$  3D views are sets of  $m$  points denoted  $\mathbf{Q}_{tj}$ , where  $t$  is the time index and  $j$  the point index. We do not use homogeneous coordinates, *e.g.*  $\mathbf{Q}_{tj}$  is a 3-

vector. The identity matrix of size  $(s \times s)$  is written  $I_{(s)}$ , the zero matrix  $0$  and the zero vector  $\mathbf{0}$ . We use  $I$  for the  $(3 \times 3)$  identity matrix. The Kronecker product is written  $\otimes$ , matrix Frobenius norm as  $\|\cdot\|$  and the Moore-Penrose pseudo-inverse as  $^\dagger$ .

### 3. Non-Rigid Shape and Pose

#### 3.1. Non-Rigid Shape

We describe the low-rank non-rigid shape model. The pose of the 3D sensor is modeled by 3D Euclidean transformations  $\{(R_t, \mathbf{y}_t)\}$  with  $R_t$  an orthonormal matrix and  $\mathbf{y}_t \in \mathbb{R}^3$  such that  $\hat{\mathbf{Q}}_{tj} = R_t \tilde{\mathbf{Q}}_{tj} + \mathbf{y}_t$ . The  $\{\tilde{\mathbf{Q}}_{tj}\}$  form a motionless version of the 3D views, *i.e.* that do not undergo any ‘global motion’, but are deforming through time. The low-rank shape model represents the  $\{\tilde{\mathbf{Q}}_{tj}\}$  as linear combinations of  $l$  basis shapes  $\{\mathbf{B}_{kj}\}$ :  $\tilde{\mathbf{Q}}_{tj} = \sum_{k=1}^l \xi_{tk} \mathbf{B}_{kj}$ . The time-varying  $\{\xi_{tk}\}$  are the *configuration weights*. Introducing the  $\{(R_t, \mathbf{y}_t)\}$ , we obtain the *explicit model*:

$$\hat{\mathbf{Q}}_{tj} = R_t \left( \sum_{k=1}^l \xi_{tk} \mathbf{B}_{kj} \right) + \mathbf{y}_t \quad (1)$$

$$= M_t \mathbf{B}_j + \mathbf{y}_t \quad \text{with} \quad (2)$$

$$M_t = R_t \begin{pmatrix} \xi_{t1} I & \cdots & \xi_{tl} I \end{pmatrix}. \quad (3)$$

We call  $M_t$  a  $(3 \times r)$  *explicit non-rigid motion matrix* and  $\mathbf{B}_j = (\mathbf{B}_{1j}^\top \cdots \mathbf{B}_{lj}^\top)$  a  $(r \times 1)$  *non-rigid basis shape vector*. Parameter  $r = 3l$  is the *rank* of the model. For reasons that are made clearer below, we derive a bilinear *implicit model*. Let  $\mathcal{A}$  be a  $(3l \times 3l)$  rank- $3l$  matrix. It is seen that  $\hat{\mathbf{Q}}_{tj} = M_t \mathbf{B}_j + \mathbf{y}_t = (M_t \mathcal{A}^{-1})(\mathcal{A} \mathbf{B}_j) + \mathbf{y}_t$ , yielding:

$$\hat{\mathbf{Q}}_{tj} = N_t \mathbf{S}_j + \mathbf{y}_t, \quad (4)$$

with  $N_t = M_t \mathcal{A}^{-1}$  and  $\mathbf{S}_j = \mathcal{A} \mathbf{B}_j$ . We call  $N_t$  and  $\mathbf{S}_j$  the *implicit non-rigid motion matrix* and *shape vector*, and  $\mathcal{A}$  a *corrective transformation matrix*.

#### 3.2. Non-Rigid Pose

Pose in a non-rigid environment has a rigid and a non-rigid counterpart. The rigid part  $\{(R_t, \mathbf{y}_t)\}$  represents the ‘global’ motion of the environment relative to the sensor. It gives the ‘true’ *relative* sensor displacement. In contrast, the non-rigid part only concerns the environment, and not the imaging sensor. In the above-described model, it is represented by the configuration weights  $\{\xi_{tk}\}$ , giving the intrinsic, *i.e.* motionless, deformations of the environment at each time instant. The motionless and deformationless environment is modeled by the basis shapes  $\{\mathbf{B}_{kj}\}$ .

The implicit model is useless for pose estimation: it can be seen as an ‘uncalibrated’ model of the environment. However, its ML (Maximum Likelihood) Estimate can be computed very reliably, as will be seen in the next section.

## 4. Learning the Environment

Given a collection of 3D views, we learn the environment by estimating the parameters of the low-rank shape model. Note that only the basis shapes  $\{\mathbf{B}_{kj}\}$  are subsequently used for pose estimation, see §5. However, to get an ML Estimate, all parameters of the model must be computed.

We state the ML residual error and show how to compute the translations. We first tackle the case of the implicit model and then the explicit one. We assume all points to be visible in all 3D views.

#### 4.1. Maximum Likelihood residual error

Assuming that the error on the 3D points is Gaussian, centred and i.i.d., the ML residual error is:

$$\mathcal{D}^2 = \frac{1}{nm} \sum_{t=1}^n \sum_{j=1}^m d^2(\hat{\mathbf{Q}}_{tj}, \mathbf{Q}_{tj}), \quad (5)$$

where  $d^2(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\|^2$  is the Euclidean distance measure and  $\{\hat{\mathbf{Q}}_{tj}\}$  are corrected points, exactly explained by the non-rigid shape model.

#### 4.2. Computing the Translations

We show that the translations  $\mathbf{y}_t$  can be eliminated prior to estimating the other parameters. By substituting equation (1) or equation (4) in the residual error (5) and nullifying its partial derivatives with respect to  $\mathbf{y}_t$ , we obtain  $\mathbf{y}_t = \frac{1}{m} \left( \sum_{j=1}^m \mathbf{Q}_{tj} - \hat{\mathbf{Q}}_{tj} \right)$ . The origin of the  $r$ -dimensional space containing the non-rigid shape vectors is arbitrary and is chosen such that  $\sum_{j=1}^m \mathbf{S}_j = \mathbf{0}$  in the implicit case and  $\sum_{j=1}^m \mathbf{B}_j = \mathbf{0}$  in the explicit case, giving for the translation  $\mathbf{y}_t$  the centroid  $\mathbf{y}_t = \frac{1}{m} \sum_{j=1}^m \mathbf{Q}_{tj} = \bar{\mathbf{Q}}_t$  of the  $t$ -th 3D view. This means that one cancels the translations out by centring each set of 3D points on its centroid:  $\mathbf{Q}_{tj} \leftarrow \mathbf{Q}_{tj} - \bar{\mathbf{Q}}_t$ . Henceforth, we assume that this has been done.

#### 4.3. Shape Learning With the Implicit Model

We consider the implicit non-rigid shape model of equation (4). We factorize the 3D views  $\{\mathbf{Q}_{tj}\}$  into implicit non-rigid motion matrices  $\{N_t\}$  and shape vectors  $\{\mathbf{S}_j\}$ . The problem is to minimize the ML residual error (5) over the  $\{\hat{\mathbf{Q}}_{tj}\}$  such that  $\hat{\mathbf{Q}}_{tj} = N_t \mathbf{S}_j$ . Rewrite (5) as:

$$\mathcal{D}^2 \propto \|\hat{\mathcal{Q}} - \mathcal{Q}\|^2,$$

where  $\mathcal{Q}$  is the  $(3n \times m)$  *measurement matrix*:

$$\mathcal{Q} = \begin{pmatrix} \mathbf{Q}_{11} & \cdots & \mathbf{Q}_{1m} \\ \vdots & \ddots & \vdots \\ \mathbf{Q}_{n1} & \cdots & \mathbf{Q}_{nm} \end{pmatrix},$$

and  $\hat{Q}$  is defined by the implicit  $(3n \times 3l)$  ‘non-rigid joint motion matrix’  $\mathcal{N}$  and the  $(3l \times m)$  ‘non-rigid joint structure matrix’  $\mathcal{S}$  as  $\hat{Q} = \mathcal{N}\mathcal{S}$  with  $\mathcal{N}^\top = (\mathbf{N}_1^\top \cdots \mathbf{N}_n^\top)$  and  $\mathcal{S} = (\mathbf{S}_1 \cdots \mathbf{S}_m)$ . Since  $\mathcal{N}$  has  $3l$  columns and  $\mathcal{S}$  has  $3l$  rows,  $\hat{Q}$  has maximum rank  $3l$ . The problem is to find the closest rank- $3l$  matrix  $\hat{Q}$  to  $Q$ . Let  $Q = U\Sigma V^\top$  be a Singular Value Decomposition (SVD) of matrix  $Q$ , see *e.g.* [6], where  $U$  and  $V$  are orthonormal matrices and  $\Sigma$  is diagonal and contains the singular values of  $Q$ . Let  $\Sigma = \Sigma_u \Sigma_v$  be any decomposition of  $\Sigma$ , *e.g.*  $\Sigma_u = \Sigma_v = \sqrt{\Sigma}$ . The non-rigid joint motion and structure matrices are obtained by, loosely speaking, ‘truncating’ the decomposition by nullifying all but the  $3l$  largest singular values, which leads, assuming the singular values in decreasing order in  $\Sigma$ , to  $\mathcal{N} = \psi_{3l}(U\Sigma_u)$  and  $\mathcal{S} = \psi_{3l}^\top(V\Sigma_v^\top)$ , where  $\psi_c(W)$  is formed with the  $c$  leading columns of matrix  $W$ .

#### 4.4. Shape Learning With the Explicit Model

The aim is to compute the ML Estimate of the configuration weights, rotation matrices and non-rigid structure in equation (1) by minimizing the residual error (5). This is a nonlinear problem for which two approaches have been followed in the non-rigid factorization literature. On the one hand Bregler *et al.* [4], Brand [3], Aanaes *et al.* [1], Del Bue *et al.* [5] and Xiao *et al.* [13] compute a matrix  $\mathcal{A}$  that upgrades the implicit motion matrix  $\mathcal{N}$  so that the metric constraints of the explicit model are enforced. Xiao *et al.* show that in order to get the correct solution, two types of metric constraints must be taken into account: the *rotation constraints* and the *basis constraints*, from which they derive a closed-form solution for matrix  $\mathcal{A}$ .

On the other hand, Torresani *et al.* [11] directly learn the parameters of the explicit model. They propose a comprehensive system based on a generalized EM (Expectation Maximization) algorithm. An important, still unsolved problem is to find a suitable initialization, since EM performs local optimization only.

Our solution lies in the second category: a suboptimal initialization is computed and subsequently refined in a bundle adjustment manner. These two steps are presented below, followed by an analysis of the ambiguities of the solution.

##### 4.4.1. Initializing

**The rotations.** Brand proposes a solution based on upgrading the implicit motion matrices [3], which requires at least  $n \geq \frac{l(9l+3)}{4}$  3D views to compute a corrective transformation and is thus not feasible for many practical cases. For example, at least 39 views giving independent constraints are necessary to use this method with the sequence presented in §6.2. In [5], the authors compute a block-diagonal corrective transformation matrix. Another

solution used in [1] is to assume that the environment has a sufficiently strong rigid component, and to estimate the rotation using a standard procedure such as [8]. This approach is not feasible for highly deforming environments.

In contrast, we propose an approach taking the non-rigid nature of the environment into account. Our algorithm is presented below in the occlusion-free case for simplicity, but can be easily extended to the missing data case. Consider the explicit non-rigid joint motion equation  $Q = \mathcal{M}\mathcal{B}$  with:

$$\mathcal{M} = \begin{pmatrix} \xi_{11}R_1 & \cdots & \xi_{1l}R_l \\ \vdots & \ddots & \vdots \\ \xi_{n1}R_n & \cdots & \xi_{nl}R_n \end{pmatrix} \quad \text{and} \quad \mathcal{B} = \begin{pmatrix} \mathbf{B}_{11} & \cdots & \mathbf{B}_{1m} \\ \vdots & \ddots & \vdots \\ \mathbf{B}_{l1} & \cdots & \mathbf{B}_{lm} \end{pmatrix}.$$

Define two subsets  $\mathbb{A}$  and  $\mathbb{B}$  of  $n_a$  and  $n_b$  3D views respectively,  $Q_a = \mathcal{M}_a\mathcal{B}$  and  $Q_b = \mathcal{M}_b\mathcal{B}$ . Our goal is to eliminate the structure  $\mathcal{B}$  from the equations. We assume without loss of generality  $\text{rank}(Q_b) \geq 3l$ . This implies  $n_b \geq l$ . We express  $\mathcal{B}$  in terms of  $Q_b$  and  $\mathcal{M}_b$  using the equation subset  $\mathbb{B}$  as  $\mathcal{B} = \mathcal{M}_b^\dagger Q_b$ . Plugging this into the equation subset  $\mathbb{A}$  yields  $Q_a = \mathcal{M}_a\mathcal{B} = \mathcal{M}_a\mathcal{M}_b^\dagger Q_b$  that we rewrite:

$$\underbrace{(\mathbf{I}_{(3n_a)} - (\mathcal{M}_a\mathcal{M}_b^\dagger))}_{\mathcal{Z}} Q_{ab} = 0_{(3n_a \times m)} \quad (6)$$

where  $n_{ab} = n_a + n_b$  and  $Q_{ab}^\top = (Q_a^\top \ Q_b^\top)$ . We call matrix  $\mathcal{Z}_{(3n_a \times 3n_{ab})}$  a *3D non-rigid tensor*. Let us examine more closely the expression of  $\mathcal{M}_a\mathcal{M}_b^\dagger$ . The joint motion matrix can be rewritten as  $\mathcal{M} = \mathcal{R}(\Xi \otimes \mathbf{I})$  where  $\mathcal{R} = \text{diag}(R_1, \dots, R_n)$  is an orthonormal matrix and  $\Xi$  is an  $(n \times l)$  matrix containing the  $\{\xi_{ik}\}$ . Similarly,  $\mathcal{M}_a = \mathcal{R}_a(\Xi_a \otimes \mathbf{I})$  and  $\mathcal{M}_b = \mathcal{R}_b(\Xi_b \otimes \mathbf{I})$ , yielding:

$$\begin{aligned} \mathcal{M}_a\mathcal{M}_b^\dagger &= \mathcal{R}_a(\Xi_a \otimes \mathbf{I})(\mathcal{R}_b(\Xi_b \otimes \mathbf{I}))^\dagger \\ &= \mathcal{R}_a(\Xi_a \otimes \mathbf{I})(\Xi_b \otimes \mathbf{I})^\dagger \mathcal{R}_b^\top, \end{aligned}$$

since  $\mathcal{R}_b$  is an orthonormal matrix. We make use of the following properties:  $(S \otimes \mathbf{I})^\dagger = S^\dagger \otimes \mathbf{I}$  and  $(S \otimes \mathbf{I})(S' \otimes \mathbf{I}) = (SS') \otimes \mathbf{I}$  to get:

$$\mathcal{M}_a\mathcal{M}_b^\dagger = \mathcal{R}_a \left( (\Xi_a \Xi_b^\dagger) \otimes \mathbf{I} \right) \mathcal{R}_b^\top.$$

Substituting in equation (6) and multiplying on the left by the orthonormal  $\mathcal{R}_a^\top$  yields:

$$\left( \mathcal{R}_a^\top - \left( (\Xi_a \Xi_b^\dagger) \otimes \mathbf{I} \right) \mathcal{R}_b^\top \right) Q_{ab} = 0_{(3n_a \times m)}. \quad (7)$$

From this equation, knowing  $\mathcal{R}_a$  and using the orthonormality constraints on  $\mathcal{R}_b$  to eliminate the weights  $\Xi_a \Xi_b^\dagger$  should allow to compute  $\mathcal{R}_b$ . We use the fact that the coordinate frame can be aligned on a reference view  $t$ , *i.e.* such that  $R_t = \mathbf{I}$  and choose one view in the initial set of 3D views  $\mathbb{A}$  to be the reference one.

The first idea that comes to mind to solve this problem is to consider the left nullspace of  $\mathcal{Q}_{ab}$ . Define a  $(3n_{ab} \times (3n_{ab} - 3l))$  matrix  $\mathbf{U}$  whose columns span the left nullspace of  $\mathcal{Q}_{ab}$ :  $\mathbf{U}^\top \mathcal{Q}_{ab} = 0$ . Using equation (7), we obtain:

$$\left( \mathcal{R}_a^\top - \left( (\Xi_a \Xi_b^\top) \otimes \mathbf{I} \right) \mathcal{R}_b^\top \right) = \mathbf{H} \mathbf{U}^\top,$$

where  $\mathbf{H}$  accounts for the fact that any linear combination of the columns of  $\mathbf{U}$  are in the left nullspace of  $\mathcal{Q}_{ab}$ . While this approach works fine in the absence of noise contaminating the data, it is however very unstable and useless when even very slight noise is present in the data. Indeed, if one employs *e.g.* SVD to compute matrix  $\mathbf{U}$ , then the singular vectors corresponding to the lowest singular values will be selected, and will not in general allow to recover the sought-after rotations, since the SVD mixes the singular vectors to obtain the lowest residual error as possible.

The second idea that comes to mind is to estimate each rotation in  $\mathbb{B}$  and the corresponding weight at a time. Consider a 3D view  $g \in \mathbb{A}$ . Equation (7) induces the following residual error:

$$\sum_{j=1}^m \|\mathbf{R}_g^\top \mathbf{Q}_{gj} - \sum_{t \in \mathbb{B}} \zeta_t \mathbf{R}_t^\top \mathbf{Q}_{tj}\|^2, \quad (8)$$

where  $\{\zeta_t\}$  are unknown weights. Initialize all rotations in  $\mathcal{R}_b$  to the identity:  $\mathbf{R}_t^0 = \mathbf{I}$ ,  $t \in \mathbb{B}$ . Let  $p \leftarrow 0$  be the iteration counter. The idea is to iteratively compute the  $t$ -th rotation for  $t \in \mathbb{B}$  while holding the other  $n_b - 1$  rotations in  $\mathbb{B}$  until convergence, by minimizing the residual error (8) that we rewrite:

$$\sum_{j=1}^m \|\mathbf{E}_j^p - \zeta_t^{p+1} (\mathbf{R}_t^{p+1})^\top \mathbf{Q}_{tj}\|^2 \quad (9)$$

with:

$$\mathbf{E}_j^p = \mathbf{R}_g^\top \mathbf{Q}_{gj} - \left( \sum_{t \in \mathbb{B}} (\mathbf{S}_t^p)^\top \mathbf{Q}_{tj} \right), \quad (10)$$

where  $\mathbf{S}_t^p$  is the latest estimate, *i.e.* :

$$\mathbf{S}_t^p = \begin{cases} \zeta_t^{p+1} \mathbf{R}_t^{p+1} & \text{if it is computed} \\ \zeta_t^p \mathbf{R}_t^p & \text{otherwise.} \end{cases}$$

We use a standard procedure for computing the 3D rotation and scale from 3D point correspondences – here  $\{\mathbf{E}_{ij}^p \leftrightarrow \mathbf{Q}_{ij}\}$  – due to [8] to solve this problem. Our algorithm is summarized in table 1. Note that at most  $l$  rotations in  $\mathcal{R}_b$  can be computed at each iteration which implies that the number of rotations in  $\mathcal{R}_b$  must be  $l$ . This is why only the smallest, *i.e.* minimal 3D non-rigid tensors can be used by our algorithm. Also, the unknown  $\mathcal{M}_b$  must be full-rank. We use the corresponding implicit  $\mathcal{N}_b$  to check that this is the case, since there exists a full-rank corrective transformation matrix  $\mathcal{A}$  such that  $\mathcal{N}_a \mathcal{A} = \mathcal{M}_a$ . In the case of missing data, the sum in equations (8) and (9) is simply replaced by a sum over the points seen in subsets  $\mathbb{A}$  and  $\mathbb{B}$ .

## OBJECTIVE

Given  $n$  3D views  $\{\mathbf{Q}_{tj}\}$  of  $m$  corresponding points and the rank  $3l$  of the non-rigid model, compute the relative pose  $\{(\mathbf{R}_t, \mathbf{y}_t)\}$  of the 3D sensor, the non-rigid pose of the environment, *i.e.* the configuration weights  $\{\xi_{tk}\}$ , while learning the low-rank non-rigid shape model  $\{\mathbf{B}_{kj}\}$ .

## ALGORITHM

1. **Set initial equation sets.**  $\mathbb{A}$  is any 3D view  $t$ ,  $\mathbb{B}$  is any  $l$  3D views at least one of them not in  $\mathbb{A}$  and such that  $\mathcal{N}_b$  is full-rank,  $\mathbf{R}_t \leftarrow \mathbf{I}$  and  $\mathcal{R}_a \leftarrow \mathbf{I}$ .
2. **Compute the rotations:**
  - (a) Set initial rotations  $\mathbf{R}_t^0 \leftarrow \mathbf{I}$ ,  $t \in \mathbb{B}$  and the iteration counter  $p \leftarrow 0$ .
  - (b) For  $t \in \mathbb{B}$ : form the  $\{\mathbf{E}_j^p\}$ , equation (10). Compute  $\mathbf{R}_t^{p+1}$  by minimizing (9), see Horn *et al.* [8].
  - (c)  $p \leftarrow p + 1$ .
  - (d) If the decrease in the residual error is smaller than  $\varepsilon$ , go to step 3 else go to step b.
3. **Test convergence.** If all rotations are computed, stop.
4. **Update equation sets.**  $\mathbb{A} \leftarrow \mathbb{A} \cup \mathbb{B}$  and  $\mathbb{B}$  is any  $l$  3D views, at least one of them not in  $\mathbb{A}$  and such that  $\mathcal{N}_b$  is full-rank.
5. **Iterate.** go to step 2.

Table 1: The proposed initialization algorithm for the explicit model parameters.

**The configuration weights and non-rigid structure.** Consider the ML residual error (5) that we rewrite below for convenience:

$$\mathcal{D}^2 = \frac{1}{nm} \sum_{t=1}^n \sum_{j=1}^m \|\mathbf{Q}_{tj} - \mathbf{R}_t \left( \sum_{k=1}^l \xi_{tk} \mathbf{B}_{kj} \right)\|^2.$$

Let  $\tilde{\mathbf{Q}}_{tj} = \mathbf{R}_t^\top \mathbf{Q}_{tj}$  be a motionless version of the 3D points, the residual error transforms in:

$$\mathcal{D}^2 = \frac{1}{nm} \sum_{t=1}^n \sum_{j=1}^m \|\tilde{\mathbf{Q}}_{tj} - \left( \sum_{k=1}^l \xi_{tk} \mathbf{B}_{kj} \right)\|^2.$$

Introduce matrices  $\mathcal{L}_{(n \times 3m)}$  and  $\mathcal{T}_{(l \times 3m)}$  which are obtained by reorganizing  $\tilde{\mathbf{Q}}$  and  $\mathcal{B}$ , respectively:

$$\tilde{\mathcal{L}} = \begin{pmatrix} \tilde{\mathbf{Q}}_{11}^\top & \cdots & \tilde{\mathbf{Q}}_{1m}^\top \\ \vdots & \ddots & \vdots \\ \tilde{\mathbf{Q}}_{n1}^\top & \cdots & \tilde{\mathbf{Q}}_{nm}^\top \end{pmatrix} \quad \text{and} \quad \mathcal{T} = \begin{pmatrix} \mathbf{B}_{11}^\top & \cdots & \mathbf{B}_{1m}^\top \\ \vdots & \ddots & \vdots \\ \mathbf{B}_{l1}^\top & \cdots & \mathbf{B}_{lm}^\top \end{pmatrix}.$$

The residual error is rewritten  $\mathcal{D}^2 = \frac{1}{nm} \|\tilde{\mathcal{L}} - \Xi \mathcal{T}\|^2$ . This means that matrix  $\tilde{\mathcal{L}}$  has rank  $l$  at most. Similarly to §4.3, let  $\tilde{\mathcal{L}} = \mathbf{U} \Sigma \mathbf{V}^\top$  be an SVD of matrix  $\tilde{\mathcal{L}}$ , we get  $\Xi = \psi_l(\mathbf{U} \Sigma_u)$  and  $\mathcal{T} = \psi_l^\top(\mathbf{V} \Sigma_v^\top)$ .

#### 4.4.2. Bundle Adjustment

Starting from the above-derived initial solution, we minimize the ML residual error (5) using nonlinear least-squares in a bundle adjustment manner, see *e.g.* [12]. We use the Levenberg-Marquardt algorithm, implemented to exploit the sparse block structure of the Jacobian and (Gauss-Newton approximation of) the Hessian matrices. The former is illustrated for a simple case on figure 1. Bundle ad-

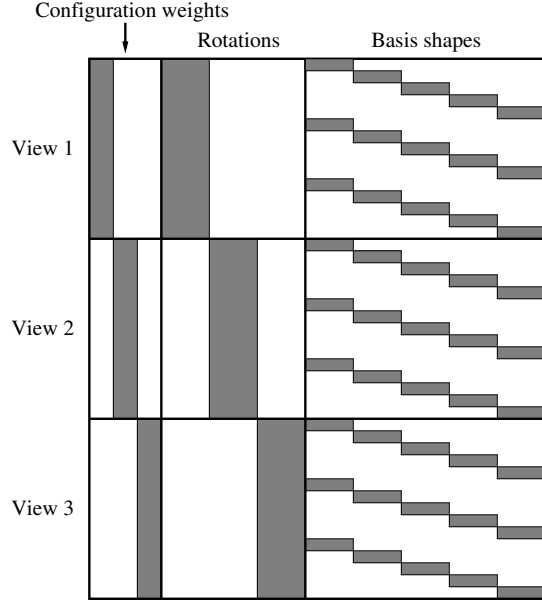


Figure 1: The Jacobian has a sparse block structure shown here for  $l = 2$  basis shapes,  $n = 3$  views and  $m = 5$  points.

justment in the non-rigid case is developed in [1, 5], where the authors show that compared to the rigid case, additional ‘gauge freedoms’ in the recovered structure and motion must be handled. However, the Levenberg-Marquardt optimization engine deals with those by damping the approximated Hessian matrix which makes it full rank. We found that the regularization term employed in [5] does not have a significant effect on the results we obtained. This is mainly due to the fact that we directly use 3D data, while [5] use image points.

#### 4.4.3. Ambiguities of the Solution

The ambiguity of the solution demonstrated by Xiao *et al.* [13] in the 2D case when only the rotation constraints are used does not hold for our algorithm. The reason is that it enforces the replicated block structure of the joint motion matrix  $\mathcal{M}$ , which provides stronger constraints than the rotation constraints only. The ambiguity matrix  $\mathcal{E}$  on the learnt model is  $\mathcal{E} = \text{diag}_l(\mathbf{R})(\Lambda_{(l \times l)} \otimes \mathbf{I})$ , where  $\text{diag}_l(\mathbf{R})$  is a  $l$  block diagonal matrix for some 3D rotation matrix  $\mathbf{R}$ ,

representing the indeterminateness of the orientation for the global coordinate frame. Matrix  $\Lambda_{(l \times l)} \otimes \mathbf{I}$  models linear combinations of the basis shapes. This shows that it is not possible to recover the ‘true’ basis shapes and configuration weights, but that ‘true’ camera pose can still be computed. The proof of this result is omitted due to lack of space.

### 5. Computing Pose

Given the non-rigid model of the environment – the basis shapes  $\{\mathbf{B}_{kj}\}$  – and a 3D view  $\{\mathbf{Q}_j\}$ , we want to estimate the pose of the 3D sensor, namely the Euclidean transformation  $(\mathbf{R}, \mathbf{y})$ , jointly with the non-rigid counterpart of the pose, *i.e.* the configuration weights  $\{\xi_k\}$ . Note that we drop index  $t$  since only one 3D view is considered in this section. It is not necessary to observe all points used in the learning phase to compute pose. The ML residual error is:

$$\mathcal{C}^2 = \frac{1}{m} \sum_{j=1}^m d^2(\mathbf{M}\mathbf{B}_j + \mathbf{y}, \mathbf{Q}_j). \quad (11)$$

It must be minimized over  $(\mathbf{R}, \mathbf{y})$  and  $\{\xi_k\}$ . Matrix  $\mathbf{M}$  is defined by equation (3).

We propose to nonlinearly minimize the ML residual error (11) using the Levenberg-Marquardt algorithm. It is not possible to use a direct estimator as [8] due to the configuration weights. Note that, as shown below, the translation  $\mathbf{y}$  can be eliminated from the equation. The minimization is thus performed over  $\mathbf{R}$  and  $\{\xi_k\}$ . Such an algorithm as Levenberg-Marquardt requires one to provide an initial solution. Our algorithm for finding it is described below.

**Eliminating the translation.** The derivatives of the ML residual error (11) with respect to  $\mathbf{y}$  must vanish:  $\frac{\partial \mathcal{C}^2}{\partial \mathbf{y}} = \mathbf{0}$ , which leads to  $\mathbf{y} = \frac{1}{m} \sum_{j=1}^m (\mathbf{Q}_j - \mathbf{M}\mathbf{B}_j)$ . This result means that  $\mathbf{y}$  is given by the difference between the centroid of the points  $\{\mathbf{Q}_j\}$  and the centroid predicted by the points from the shape model  $\{\mathbf{M}\mathbf{B}_j\}$ , which vanishes if the set of points used for computing the pose is exactly the same as the one used in the learning phase. In any case, by centering the points on their centroid, the translation vanishes. Henceforth, we assume this has been done and rewrite the ML residual error (11) as:

$$\mathcal{C}^2 = \frac{1}{m} \sum_{j=1}^m d^2(\mathbf{M}\mathbf{B}_j, \mathbf{Q}_j). \quad (12)$$

**Initializing the rotation and configuration weights.** We linearly compute a motion matrix  $\tilde{\mathbf{M}}$  without enforcing the correct replicated structure by  $\min_{\tilde{\mathbf{M}}} \sum_{j=1}^m d^2(\tilde{\mathbf{M}}\mathbf{B}_j, \mathbf{Q}_j)$ , which yield:

$$\tilde{\mathbf{M}} = (\mathbf{Q}_1 \cdots \mathbf{Q}_m) (\mathbf{B}_1 \cdots \mathbf{B}_m)^\dagger.$$

We extract the  $\{\xi_k\}$  and  $R$  from  $\tilde{M}$  by solving  $\min_{R, \{\xi_k\}} \sum_{k=1}^l \|\tilde{M}_k - \xi_k R\|^2$ , where the  $\tilde{M}_k$  are  $(3 \times 3)$  blocks from  $\tilde{M}$ . By vectorizing and reorganizing the residual error, we obtain:

$$\left\| \underbrace{\begin{pmatrix} \text{vect}^T(\tilde{M}_1) \\ \vdots \\ \text{vect}^T(\tilde{M}_l) \end{pmatrix}}_{\Lambda} - \underbrace{\begin{pmatrix} \xi_1 \\ \vdots \\ \xi_l \end{pmatrix}}_{\xi} \underbrace{\text{vect}^T(\tilde{R})}_{\tilde{r}} \right\|^2,$$

which is a rank-1 approximation problem that we solve by ‘truncating’ the SVD  $\Lambda = U\Sigma V^T$ , as in §4.3:  $\xi = \psi_1(U\Sigma)$  and  $\tilde{r} = \psi_1^T(V)$ . Note that  $\|\tilde{r}\| = \|\tilde{R}\| = 1$ . Matrix  $\tilde{R}$  must be subsequently corrected to give  $R$  by enforcing the orthonormality constraints. This is done by finding the closest orthonormal matrix to  $\tilde{R}$  using SVD, see [8]:  $\tilde{R} = U\Sigma V^T$  gives  $R = \frac{1}{3} \text{tr}(\Sigma) \det(U) \det(V) UV^T$ , while compensating the possible sign change by  $\xi \leftarrow \det(U) \det(V) \xi$ .

## 6. Experimental Evaluation

### 6.1. Simulated Data

We report experimental results on simulated data. The default simulation setup consists of  $n = 15$  time-varying 3D views, each containing  $m = 35$  points. They are generated by randomly drawn linear combinations of  $l = 3$  basis shapes, all of them lying in a sphere with unit radius. An additive, zero-mean Gaussian noise with variance  $\sigma = 0.01$  (*i.e.* 1% of the scene scale) is added to the 3D points. We vary each of these parameters in turn. We average the error measures over 100 trials. The true number of basis shapes is used by the algorithms.

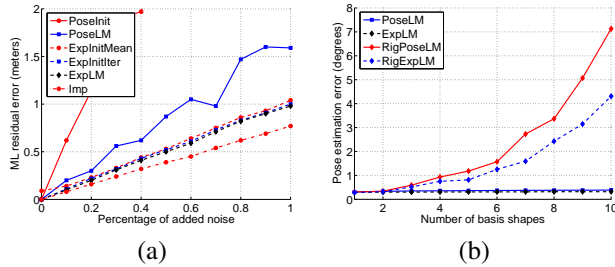


Figure 2: (a) ML residual error against the level of added noise and (b) pose error against the number of basis shapes.

**Environment learning.** We observe on figure 2 (a) that the ML residual error is very close to  $\sigma$ . The implicit learning IMP consistently gives a significantly lower residual error than the explicit learning algorithms  $\text{EXP}\star$ . This means that, despite the fact that the data were generated using the explicit low-rank shape model, the extra degrees of freedom

of the implicit model represent quite well the added Gaussian noise.

We observe that the difference between the three explicit learning methods is small compared to the difference with IMP. EXPLM (from §4.4.2, ‘LM’ stands for Levenberg-Marquardt) always performs better than EXPINITITER (from table 1), which always performs better than EXPINITMEAN (based on [8] to get the rotations). This means that the residual error (8), which is minimized by EXPINITITER while estimating the minimal 3D non-rigid tensors, is well-adapted to our problem.

Figure 2 (b) compares the error raised by the rotation part of the pose, in degrees, between our non-rigid algorithms and rigid SFM and pose algorithms, respectively dubbed RIGEXPLM and RIGPOSELM. We observe that the proposed EXPLM gives errors independent of the number of basis shapes, while, as could have been expected, RIGEXPLM rapidly degrades as the number of basis shapes grows.

**Pose computation.** Figure 2 (a) shows that all pose algorithms  $\text{POSE}\star$  consistently give a higher residual error than the explicit learning algorithms. This is explained by the fact that pose estimation suffers from the errors in the learnt model *and* in the 3D view.  $\text{POSEINIT}$  gives quite high errors, roughly  $5\sigma$ , while  $\text{POSELM}$  converges to roughly  $1.5\sigma$  which is reasonable. The same remarks as for the learning algorithms can be made in the case of pose, for figure 2 (b).

Another experiment was intended to assess to which extent, reliable pose estimate can be obtained when the environment is deforming in a very different way compared to the learning stage. Let  $\nu$  be the mean value of the configuration weights. We alter them by adding randomly drawn perturbations with increasing magnitude  $\mu$ , and generate a 3D view with these parameters, from which pose is estimated. Obviously, the results depend on the simulation setup, the number of points, views, basis shapes and the level of noise. However, we observe that for  $\mu \leq 1.3\nu$ , the residual error indicates that the pose estimate is reasonable for most configurations. For  $\mu > 1.3\nu$ , the pose estimate rapidly degrades.

### 6.2. Real Data

We tested our algorithms on sets of 3D points reconstructed from a calibrated stereo rig. The sequence consists of  $n = 650$  pairs of views. The  $m = 30$  point tracks were obtained semi-automatically and reconstruction was performed using ML triangulation, *i.e.* by minimizing the reprojection error. The reprojection error we obtained is 4.7276 pixels, which is rather large and explained by the low quality of the manually entered point tracks.

We used a subset of the full sequence, made of 1 3D view over 25 from 1 to 551, that is 23 3D views, for learning the environment. The remaining 3D views are registered by

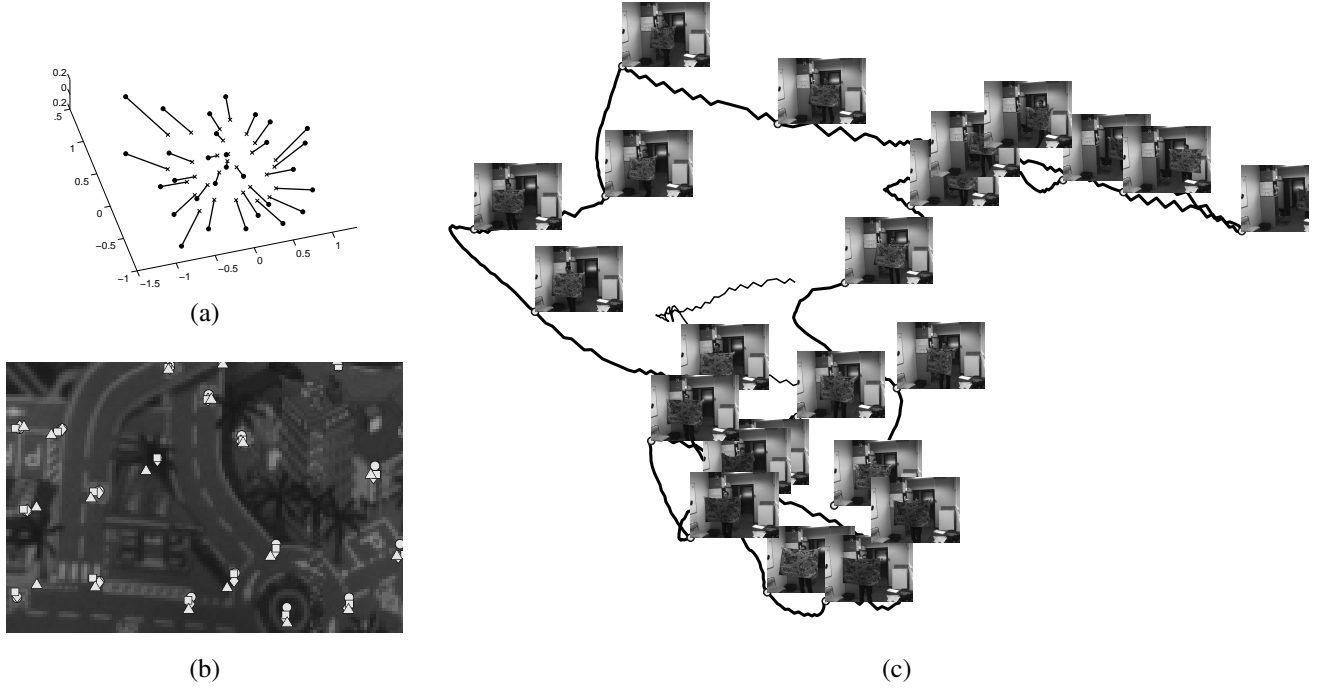


Figure 3: (a) The first (circles) and the second (crosses) basis shapes. (b) Zoom on one of the original images: the circles are the tracks, the diamonds are the reprojections from pairwise Structure-From-Motion, the squares are the reprojections from our non-rigid shape model, and the triangles are the reprojections from a rigid shape model. (c) The path representing the translation part of the pose, together with the images corresponding to the 3D views used for learning the environment, indicated by white dots at the corner of the images. The ‘interpolated’ positions are shown using a bold line.

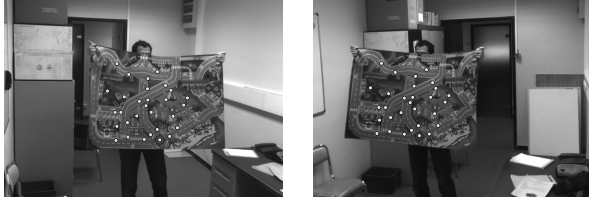


Figure 4: One out of the 650 stereo pairs used in the experiments, overlaid with the 30 point tracks.

computing pose. For views  $1 < i < 551$ , this can be viewed as ‘interpolation’ since the surrounding 3D views are used for learning the environment, while for views  $551 < i < 650$  this can be viewed as an ‘extrapolation’ of the model since new pose and deformations are seen in these views.

An important aspect is the choice of the number  $l$  of basis shapes. If  $l$  is too low, the model is not able to represent all the possible deformations, while if  $l$  is too high, the noise is modeled, resulting in unreliable pose estimates in both cases. We propose to manually choose  $l$  by examining the graphs shown on figure 5. It shows the ML residual errors and the reprojection errors, *i.e.* the Sum of Squared Differences between measured and predicted image points, resulting of the learning algorithms for different numbers

of basis shapes. We observe that the 3D ML residual er-

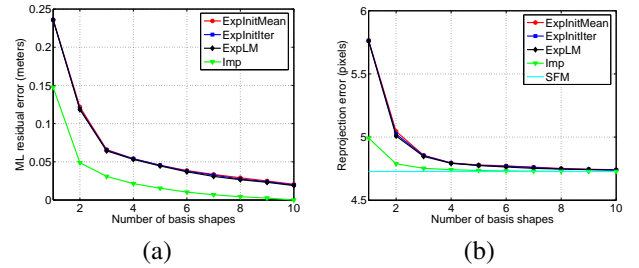


Figure 5: (a) ML residual error and (b) 2D reprojection error versus the number of basis shapes.

ror and the 2D reprojection error decrease while  $l$  increases, the former towards 0 and the latter towards the reconstruction error, shown by an horizontal line ‘SFM’ on the graph, which was expected. Based on this graph, we choose  $l = 4$ , for which the EXPLM ML residual error is 5.32 centimeters and the 2D reprojection error is 4.7922 pixels. For comparison, a rigid environment model gives a 23.58 centimeters ML residual error and a 5.7605 pixels 2D reprojection error. It is important to note that for  $l = 5$  and  $l = 6$  basis shapes, very similar pose estimates are subsequently obtained.

Figure 3 (c) shows the learnt path, together with some

key images. This result appears visually satisfying. Note that this is only intended to let the reader figure out the experimental conditions since we observed that the path, in other words the translation part of the pose  $\{\mathbf{y}_t\}$ , is very similar between the rigid and the non-rigid models, the mean difference being smaller than a centimeter. However, the mean difference in the rotations is 2.81 degrees, which is significant, but difficult to illustrate visually. The mean ML residual errors are 8.66 and 12.83 centimeters for the ‘(a) The first (circles) and the second (crosses) basis shapes. (b) Zoom on one of the original images: the circles are the tracks, the diamonds are the reprojections from pairwise Structure-From-Motion, the squares are the reprojections from our non-rigid shape model, and the triangles are the reprojections from a rigid shape model.interpolated’ and the ‘extrapolated’ poses respectively.

Figure 3 (a) shows the two first basis shapes that were learnt. We observe that the deformation is significant.

Finally, figure 3 (b) shows a zoom on points reprojected in an image from different models. We clearly see that the non-rigid shape model fits the data much better than a rigid shape model.

The computation time for the learning phase is of the order of a minute while pose estimation is roughly a tenth of a second.

## 7. Conclusions

We proposed methods for learning the parameters of a deforming environment, represented by a low-rank shape model, from data gathered by a 3D sensor. We introduced 3D non-rigid tensors, relating the positions of corresponding points in several, time-deforming, 3D views, and derived methods for estimating the parameters of the Maximum Likelihood implicit and explicit models. Based on a learnt model, we showed that the rigid pose of the 3D sensor can be computed while the environment undergoes new deformations. Experimental results on simulated data show that accurate pose estimates are obtained, even when the environment deforms in ways not observed during the learning phase. This is confirmed by results on real data, showing that the pose estimates are accurate enough to be used for tasks such as robot positioning and augmented reality.

One weakness of the approach is to rely on 3D point correspondences. We are currently working on using more robust types of inputs, such as contours or image patches, that can be reliably tracked through sequences of stereo pairs using *e.g.* particule filtering techniques. This is intended to be part of an iterative deforming environment learning system. Essential issues that will be dealt with are assessing what kind of deformations can be represented by the low-rank shape model and choosing the number of basis shapes,

which will be examined in the framework of model selection.

## References

- [1] H. Aanæs and F. Kahl. Estimation of deformable structure and motion. In *Vision and Modelling of Dynamic Scenes Workshop*, 2002.
- [2] B. Basclé and A. Blake. Separability of pose and expression in facial tracing and animation. In *International Conference on Computer Vision*, 1998.
- [3] M. Brand. Morphable 3D models from video. In *Computer Vision and Pattern Recognition*, 2001.
- [4] C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3D shape from image streams. In *Computer Vision and Pattern Recognition*, 2000.
- [5] A. Del Bue and L. Agapito. Non-rigid 3D shape recovery using stereo factorization. In *Asian Conference on Computer Vision*, 2004.
- [6] G.H. Golub and C.F. van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 1989.
- [7] J. V. Hajnal, D. L. G. Hill, and D. J. Hawkes, editors. *Medical Image Registration*. CRC Press, 2001.
- [8] B.K.P. Horn, H.M. Hilden, and S. Negahdaripour. Closed-form solution of absolute orientation using orthonormal matrices. *Journal of the Optical Society of America A*, 5(7):1127–1135, July 1988.
- [9] M. Irani. Multi-frame optical flow estimation using subspace constraints. In *International Conference on Computer Vision*, 1999.
- [10] L. Torresani and C. Bregler. Space-time tracking. In *European Conference on Computer Vision*, 2002.
- [11] L. Torresani and A. Hertzmann. Automatic non-rigid 3D modeling from video. In *European Conference on Computer Vision*, 2004.
- [12] B. Triggs, P.F. McLauchlan, R.I. Hartley, and A. Fitzgibbon. Bundle ajustment — a modern synthesis. In *Workshop on Vision Algorithms: Theory and Practice*, 2000.
- [13] J. Xiao, J.-X. Chai, and T. Kanade. A closed-form solution to non-rigid shape and motion recovery. In *European Conference on Computer Vision*, 2004.
- [14] A. J. Yezzi and S. Soatto. Deformation: Deforming motion, shape average and the joint registration and approximation of structures in images. *International Journal of Computer Vision*, 53(2):153–167, March 2003.