# Extracting Average Shapes from Occluded Non-rigid Motion

Alessio Del Bue

Institute for Systems and Robotics, Instituto Superior Tecnico, Lisbon, Portugal

**Abstract.** This paper presents a method to efficiently estimate average 3-D shapes from non-rigid motion in the case of missing data. Such a shape can be further used to accomplish full reconstruction of deformable objects and registration of non-rigid shapes. The approach is based firstly on a power method which linearly provides an initial estimate of the 3-D structure and motion components of the object shape. Secondly, non-linear optimisation is used to refine the initial linear estimation. Tests on both real and synthetic sequences show the procedure effectiveness in dealing with different degrees of occlusions in the measurements.

## 1 Introduction

Recently the inference of the 3-D structure of a deforming body viewed by an uncalibrated camera has attracted increasing interest. In a Structure from Motion (SfM) domain, non-rigid shapes have posed new problems since they violate the rigidity constraints on which previous SfM methods strongly rely. Most of the model-free approaches to non-rigid SfM available nowadays are based either on closed-form solutions [12], assuming pre-specified shape priors, or iterative non-linear optimisation techniques [5, 1, 11], requiring an appropriate initialisation in order to converge. In the latter case, average shape and motion [9] have experimentally proven to be a successful initialisation to such tasks and they can be easily computed when the full trajectory of a point lying on the deforming body is available.

However, in the case of missing data affecting the trajectories (i.e. a point being occluded for some frames) a solution for the average shape is not currently available. Estimation of structure and motion from occluded data (see [3] for a review) is an essential task for most practical applications given the difficulty to obtain complete trajectories. At this end, the solution proposed here is an iterative power method which can estimate average shapes in the case of missing data and its reliability is assessed in a full 3-D reconstruction task for deforming objects. The approach is based on the notion of average shape introduced in [9] which penalizes in a certainty-reweighted scheme the non-rigidity of trajectories. This method can be extended by reformulating power methods for SfM [7] to include the notion of non-rigidity of a trajectory and extend it to the case of missing data.

In detail, the paper firstly introduces the non-rigid factorization framework and the definition of average shape (Section 2). Then, power methods for SfM are presented in Section 3 for the case of rigidly moving objects. The new approach with missing entry is explained in Section 4 and experiments (Section 5) show its effectiveness on synthetic test and on a face modelling task.

## 2   Non-rigid Structure from Motion

### 2.1   A Factorization Approach to Deformable Modelling

Tomasi and Kanade's factorization algorithm [10] has been reformulated to the case of non-rigid 3-D structure [2]. A linear approximation of a set of $K$ basis shapes $\mathtt{S}_k$ is used to describe a 3-D time varying shape $\mathtt{X}$ such that:

$$\mathtt{X} = \sum_{k=1}^{K} l_k \mathtt{S}_k \qquad \mathtt{X}, \mathtt{S}_k \in \Re^{3 \times P} \quad l_k \in \Re \tag{1}$$

Each basis shapes $\mathtt{S}_k$ represent the mode of deformations of the deforming body and they are parameterised as a $3 \times P$ matrix which contains the 3-D locations of $P$ object points for that particular mode of deformation. Assuming an orthographic camera model the shape is then projected onto an image frame $i$ giving $P$ image points:

$$\mathtt{W}_i = \begin{bmatrix} \mathbf{w}_{i1} \ ... \ \mathbf{w}_{iP} \end{bmatrix} = \mathtt{R}_i \left( \sum_{k=1}^{K} l_{ik} \mathtt{S}_k \right) \tag{2}$$

where each $\mathbf{w}_{ij} = [u_{ij} v_{ij}]^T$ with $j = 1 \ldots P$ contains the horizontal and vertical image coordinates of the point – referred to the centroid of the object – and $\mathtt{R}_i$ encodes the first two rows of the rotation matrix for a specific frame $i$. If all $P$ points are tracked in $F$ image frames we may construct the measurement matrix $\mathtt{W}$ which can be expressed as:

$$\mathtt{W} = \begin{bmatrix} \mathbf{w}_{11} \ \ldots \ \mathbf{w}_{1P} \\ \vdots \qquad \vdots \\ \mathbf{w}_{F1} \ \ldots \ \mathbf{w}_{FP} \end{bmatrix} = \begin{bmatrix} l_{11}\mathtt{R}_1 \ \ldots \ l_{1K}\mathtt{R}_1 \\ \vdots \qquad \vdots \\ l_{F1}\mathtt{R}_F \ \ldots \ l_{FK}\mathtt{R}_F \end{bmatrix} \begin{bmatrix} \mathtt{S}_1 \\ \vdots \\ \mathtt{S}_K \end{bmatrix} = \mathtt{MS}. \tag{3}$$

Clearly, the rank of the measurement matrix is constrained to be at most $3K$, where $K$ is the number of deformations. This rank constraint can be exploited to factorize the measurement matrix into a motion matrix $\tilde{\mathtt{M}}$ and a shape matrix $\tilde{\mathtt{S}}$ by truncating the SVD of $\mathtt{W}$ to rank $3K$. However, this factorization is not unique since any invertible $3K \times 3K$ matrix $\mathtt{Q}$ can be inserted in the decomposition leading to the alternative factorization: $\mathtt{W} = (\tilde{\mathtt{M}}\mathtt{Q})(\mathtt{Q}^{-1}\tilde{\mathtt{S}})$. The focal problem to solve in non-rigid factorization schemes is to find the $\mathtt{Q}$ that renders the appropriate replicated block structure of the motion matrix and that removes the affine ambiguity, upgrading the reconstruction to a metric one.

### 2.2   Extracting Average Shapes from Deformations

Based on the framework described in the previous section, Kim & Hong [9] recently introduced a measure called the Degree of Non-rigidity (*DoN*) to estimate the deviation of a deformable point from its average position. This measure can in turn be used to extract an average shape using an iterative certainty reweighted scheme. If the average 3-D shape of a time varying shape $\mathtt{X}_i = [\mathbf{X}_{i1} \ldots \mathbf{X}_{iP}]$ is given by $\hat{\mathtt{X}} = [\hat{\mathbf{X}}_1 \ldots \hat{\mathbf{X}}_P]$ the *DoN* for point $j$ is defined as:

$$DoN_j = \sum_{i=1}^{F} (\mathbf{X}_{ij} - \hat{\mathbf{X}}_j)(\mathbf{X}_{ij} - \hat{\mathbf{X}}_j)^T. \tag{4}$$

The 2-D projection $\mathtt{C}_j$ of the $DoN$ will be thus given by:

$$\mathtt{C}_j = \sum_{i=1}^{F} \mathtt{R}_i (\mathbf{X}_{ij} - \hat{\mathbf{X}}_j)(\mathbf{X}_{ij} - \hat{\mathbf{X}}_j)^T \mathtt{R}_i^T = \sum_{i=1}^{F} (\mathbf{w}_{ij} - \hat{\mathbf{w}}_{ij})(\mathbf{w}_{ij} - \hat{\mathbf{w}}_{ij})^T \quad (5)$$

where $\mathbf{w}_{ij}$ are the image coordinates of point $j$ at frame $i$ and $\hat{\mathbf{w}}_{ij}$ are the coordinates of its projected mean shape. While the $DoN$ cannot be computed without an estimation of the mean 3-D shape, the value of its projection can be estimated directly from image measurements.

An initial estimate of the projected 2-D mean shapes $\hat{\mathbf{w}}_{ij}$ could be given simply by the first basis shape $\mathtt{S}_1$ (as in equation (3)) which could be computed with a rank-3 approximation $SVD_3(\mathtt{W}) = \hat{\mathtt{M}}\hat{\mathtt{S}}$. The projected deviation from the mean for all the points would then be defined by $\{\mathbf{w}_{ij} - \hat{\mathbf{w}}_{ij}\} = \mathtt{W} - \hat{\mathtt{M}}\hat{\mathtt{S}}$. However, a straight application of a rank-3 factorization over the first basis component does not produce an accurate measure of $\mathtt{C}_j$ as showed in [9]. To adjust the covariances, the average shape and $\mathtt{C}_j$ are iteratively estimated until convergence. However, the procedure is unable to deal with the case of missing data affecting the measurements. We will show in the next section how power methods can efficiently solve this issue.

## 3   Power Methods for Structure from Motion

SfM algorithms based on factorization require an initial decomposition of the motion $\mathtt{M}$ and structure matrix $\mathtt{S}$ given the data $\mathtt{W}$. In this context, power methods were introduced with the name *powerfactorization* by Schaffalitzky and Hartley [7] to efficiently factorise rank-constrained image measurements. This approach is an alternation method which iteratively estimates $\mathtt{M}$ and $\mathtt{S}$ by simply executing multiplications and matrices inverse. The update rules at iteration $q$ are given by [7]:

$$\mathtt{M}_q = \mathtt{W}\mathtt{S}_{q-1}^T (\mathtt{S}_{q-1}^T \mathtt{S}_{q-1})^{-1}$$

$$\mathtt{S}_q = (\mathtt{M}_q^T \mathtt{M}_q)^{-1} \mathtt{M}_q^T \mathtt{W} \quad (6)$$

They are a straightforward derivation from the orthogonal power method [6] which convergence rate depends on the ratio of the dominant singular values of $\mathtt{W}$. In the case of an affine camera viewing a moving rigid body, the update rules (6) can be modified to account for the geometrical properties of the measurements. For each frame $i = 1 \cdots F$, the projection of a point $j = 1 \cdots P$ can be expressed as:

$$\mathbf{w}_{ij} = \mathtt{A}_i \mathbf{X}_j + \mathbf{a}_i \quad (7)$$

where $\mathtt{A}_i$ is a $2 \times 3$ camera projection matrix, $\mathbf{X}_j$ a 3-vector of the 3-D coordinates and $\mathbf{a}_i$ a 2-vector of the affine camera translation. In a more compact form, equation (7) can be rewritten for every point at each frame as:

$$\mathtt{W}_i = \begin{bmatrix} \mathtt{A}_i \big| \mathbf{a}_i \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 \cdots \mathbf{X}_P \\ 1 \ \cdots \ 1 \end{bmatrix} = \mathtt{M}_i \begin{bmatrix} \mathtt{X} \\ \mathbf{1}^T \end{bmatrix} \quad (8)$$

where $\mathbf{1}$ is a $P$-vector of ones. Finally, the global expression for each frame can be written as:

$$W = \begin{bmatrix} W_1 \\ \vdots \\ W_F \end{bmatrix} = \begin{bmatrix} \left[ A_1 \middle| \mathbf{a}_1 \right] \\ \vdots \\ \left[ A_F \middle| \mathbf{a}_F \right] \end{bmatrix} \begin{bmatrix} X \\ \mathbf{1}^T \end{bmatrix} = \left[ A \middle| \mathbf{a} \right] \begin{bmatrix} X \\ \mathbf{1}^T \end{bmatrix} = MS \tag{9}$$

The algorithm for extracting the affine motion and structure of a rigid object can be summarized as follows:

> – Initialize $X_0$ with random entries.
> – Compute the $2F \times 4$ update of $M_q$ given equation (6).
> – Extract the $2F \times 1$ measurements centroid $\mathbf{a}_q$ such that $M_q = \left[ A_q \middle| \mathbf{a}_q \right]$.
> – Compute the $3 \times P$ update of $X_q$ such that: $X_q = (A_q^T A_q)^{-1} A_q^T (W - T_q)$ where $(W - T_q)$ are the centered coordinates and $T_q = \mathbf{a}_q \mathbf{1}_{1 \times P}$

## 4  Average Shape Estimation with Missing Data

### 4.1  Power Iterations and Degree of Non-rigidity

In the case of affine estimation of average shape $\hat{S}$ and motion $\hat{M}$, strongly non-rigid trajectories (which in turn provide high covariances $C_j$) are penalized in the estimation of the average components. The estimation task can be recast in the minimisation of a cost function $\chi$ such that:

$$\chi = \sum_{i,j} (\mathbf{w}_{ij} - \hat{M}_i \bar{\mathbf{X}}_j)^T C_j^{-1} (\mathbf{w}_{ij} - \hat{M}_i \bar{\mathbf{X}}_j) \tag{10}$$

where $\bar{\mathbf{X}}_j$ contains the the homogeneous coordinate for the average point such that $\bar{\mathbf{X}}_j = [\hat{\mathbf{X}}_j^T \ 1]^T$. Minimizing $\chi$ can be carried out with a minor reformulation of the power method [7]. In brief, each matrix $C_j^{-1}$ can be factored as $C_j^{-1} = B_j^T B_j$ giving:

$$\chi = \sum_{i,j} (\mathbf{w}_{ij} - \hat{M}_i \bar{\mathbf{X}}_j)^T B_j^T B_j (\mathbf{w}_{ij} - \hat{M}_i \bar{\mathbf{X}}_j) = \sum_{i,j} \| B_j \mathbf{w}_{ij} - B_j \hat{M}_i \bar{\mathbf{X}}_j \|^2 \tag{11}$$

Notice the similarity with equation (6) which hints to a solution of the minimization of (10) with a power approach. In order to obtain the two updates rules for motion and structure we can rewrite (11) such that:

$$\chi = \sum_{i,j} \| B_j \mathbf{w}_{ij} - B_j \tilde{X}_j \tilde{\mathbf{m}}_i \|^2 \tag{12}$$

with:

$$\tilde{X}_j = \begin{bmatrix} \bar{\mathbf{X}}_j^T & 0 \\ 0 & \bar{\mathbf{X}}_j^T \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{m}}_i = \begin{bmatrix} \mathbf{m}_{1i}^T & \mathbf{m}_{2i}^T \end{bmatrix}^T \tag{13}$$

where $\mathbf{m}_{1i}^T$ and $\mathbf{m}_{2i}^T$ are respectively the first and second $4 \times 1$ rows of $\hat{M}_i$. Given the quadratic costs (11) and (12) we can express the power updates for the motion as:

$$\tilde{\mathbf{m}}_i = (\sum_j \tilde{X}_j^T B_j^T B_j \tilde{X}_j)^{-1} \sum_j \tilde{X}_j^T B_j^T B_j \mathbf{w}_{ij} = (\sum_j \tilde{X}_j^T C_j^{-1} \tilde{X}_j)^{-1} \sum_j \tilde{X}_j^T C_j^{-1} \mathbf{w}_{ij} \tag{14}$$

After rearranging $\tilde{\mathbf{m}}_i \mapsto \hat{\mathsf{M}}_i = \left[ \mathsf{A}_i \middle| \mathbf{a}_i \right]$ we obtain:

$$\hat{\mathbf{X}}_j = (\sum_i \mathsf{A}_i^T \mathsf{C}_j^{-1} \mathsf{A}_i)^{-1} \sum_i \mathsf{A}_i^T \mathsf{C}_j^{-1} (\mathbf{w}_{ij} - \mathbf{a}_i) \tag{15}$$

where $\mathbf{a}_i$ is the overall translation component as defined in (9). Once the estimates for the average $\hat{\mathsf{M}}$ and $\hat{\mathsf{S}}$ are available, $\mathsf{C}_j$ is update by equation (5).

### 4.2   The Missing Data Case

We can now assume that some points are not visible in some frames due to occlusion. In order to include missing data, we can modify the power update equations in (14) and (15) to simply not include the equations regarding the missing entries giving:

$$\tilde{\mathbf{m}}_i = (\sum_j \check{\mathsf{X}}_j^T \mathsf{C}_j^{-1} \check{\mathsf{X}}_j)^{-1} \sum_j \check{\mathsf{X}}_j^T \mathsf{C}_j^{-1} \, Z_{ij} \mathbf{w}_{ij} \tag{16}$$

$$\hat{\mathbf{X}}_j = (\sum_i \mathsf{A}_i^T \mathsf{C}_j^{-1} \mathsf{A}_i)^{-1} \sum_i \mathsf{A}_i^T \mathsf{C}_j^{-1} \, Z_{ij} (\mathbf{w}_{ij} - \mathbf{a}_i) \tag{17}$$

where $Z_{ij}$ is a scalar which is zero whenever a point is missing and one otherwise. The updates have the property of efficiently estimating the centroid at each frame $\mathbf{a}_i$ since the measurement matrix of missing data may be not mean-centered. Schematically, the algorithm can be outlined as follows[1]:

---
- Initialize $\mathsf{X}$ with random entries.
- Compute the $2F \times 4$ update of $\hat{\mathsf{M}}_i$ for $i = 1 \cdots F$ given equation (16).
- Given $\hat{\mathsf{M}}_i = \left[ \mathsf{A}_i \middle| \mathbf{a}_i \right]$, extract the measurements centroid $\mathbf{a}_i$.
- Compute the update of the average 3-D structure with (17).
- Recompute $\mathsf{C}_j = \sum_{i=1}^{F} Z_{ij} (\mathbf{w}_{ij} - \mathbf{a}_i - \mathsf{A}_i \hat{\mathbf{X}}_j)(\mathbf{w}_{ij} - \mathbf{a}_i - \mathsf{A}_i \hat{\mathbf{X}}_j)^T$.
- Iterate until convergence.
---

A metric update of the average shape can be then obtained in the case of orthographic [10] and weak or para-perspective cameras [8] by computing the correct $3 \times 3$ transformation $\mathsf{Q}$ for the average shape.

### 4.3   Non-linear Optimisation and Non-rigid SfM

A full deformable 3-D reconstruction as presented in Section 2.1 can be successfully computed linearly only when particular assumptions over the data are given and when the full trajectories are available. For instance, in [12] the authors proved the existence of a unique solution and a closed form algorithm when $K$ independent 3-D shapes can be identified in the measured data. On the other hand, a more general solution consists in performing non-linear optimisation [5, 1, 11] by minimizing a cost function which reflects the full deformable model as presented in equation (3) giving:

$$\min_{\mathsf{R}_i \mathsf{S}_{kj} l_{ik}} \sum_{i,j} Z_{ij} \parallel \mathbf{w}_{ij} - \hat{\mathbf{x}}_{ij} \parallel^2 = \min_{\mathsf{R}_i \mathsf{S}_{kj} l_{ik}} \sum_{i,j} Z_{ij} \parallel \mathbf{w}_{ij} - (\mathsf{R}_i \sum_k l_{ik} \mathsf{S}_{kj}) \parallel^2 \tag{18}$$

---
[1] for clarity we drop the iteration subscript $q$

where $\mathbf{S}_{kj}$ is the $3 \times 1$ basis for the point $j$ such that $\mathbf{S}_k = [\mathbf{S}_{k1} \cdots \mathbf{S}_{kP}]$. Again, the least-squares entries for the missing data are omitted. Initialisation of the model parameters are provided by the average shape computed with our power approach.

## 5    Experiments

### 5.1    Synthetic Data

The proposed power approach is first validated using randomly generated synthetic data of a deforming shape. The 3-D bodies are generated by firstly sampling the first basis shape $\mathbf{S}_1$ over the surface of a sphere. The following basis $\mathbf{S}_2 \ldots \mathbf{S}_K$, which express the modes of deformation of the body, are generated randomly. In order to obtain a given deformation at frame $i$, the configuration weights $l_{i1} \ldots l_{iK}$ are computed by fitting 4-order polynomials to random samples, this gives more regular deformation rather then erratic motion. The computed 3-D shapes are then normalized to obtain a specific ratio of deformation defined as $\frac{\sum_{i=1}^{F} \| \sum_{k=2}^{K} l_{ik} \mathbf{S}_k \|^2}{\sum_{i=1}^{F} \| l_{i1} \mathbf{S}_1 \|^2}$ which is fixed to $0.25$. The final measurement matrix $\mathbf{W}$ is obtained by projecting each 3-D non-rigid shape onto the image plane by means of random orthographic cameras Finally, points are eliminated given different ratios of missing data.

We test the algorithm performances in providing a meaningful initialisation to the optimisation problem as defined in section 4.3. Firstly, we noticed problems in the convergence if the algorithm for computing the average shape does not include the iterative re-weighting with $\mathbf{C}_j$. The overall results are showed in table 1 with different levels of image noise affecting the data. A decrease in the algorithm's performances is given for ratios of $30\%$ and more missing data. Regarding the mean shape computation, convergence is generally achieved after 15 iterations with $10\%$ of missing data, higher ratios increase this number however, in the worst case, the algorithm was not performing more than 50 steps.

### 5.2    Real Data

The real experiments are focused on extracting a mean shape from a deforming face[2] exhibiting a light rotation and non-rigid motion especially in the mouth region We selected a 700 frames long sequence from the overall 5000 frames and 56 points are collected to form the measurement matrix $\mathbf{W}$. Occluded points appear with an overall $20\%$ ratio of missing entries. The recovered mean shape (see figure 1) is then used to initialize a full deformable reconstruction and some frames presenting the recovered 3-D depth and deformations are presented (front and side view). The approach is able to successfully recover a reasonable estimates of the depth and deformations even if the subject is not performing strong rigid motion. The final number of iterations for the power method was of 50 followed by 40 iteration of non-linear optimisation.

---

[2] sequence available at: www-prima.inrialpes.fr/FGnet/data/01-TalkingFace/talking_face.html

## 6    Conclusions

We presented a power approach to estimate average shapes from non-rigid motion in the case of missing data. Experimentally we have shown the effectiveness of the method in a deformable 3-D reconstruction task with affine cameras. The extracted average shape and motion have been shown to provide a reliable initialisation for SfM optimisation tasks in the tested cases. As a further study, the method may be extended to more general camera models (i.e. full perspective), however initial results have shown increased instability in the convergence given by the difficulty in decoupling deformations from perspective distortions. In such cases, an approach using shape priors as presented in [4] may help to successfully compute a reliable average shape.

## References

1. M. Brand. A direct method for 3d factorization of nonrigid motion observed in 2d. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, San Diego, California*, pages 122–128, 2005.
2. C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3d shape from image streams. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, Hilton Head, South Carolina*, pages 690–696, June 2000.
3. A. M. Buchanan and A. Fitzgibbon. Damped newton algorithms for matrix factorization with missing data. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, San Diego, California*, volume 2, pages 316–322, 2005.
4. A. Del Bue, X. Lladó, and L. Agapito. Non-rigid face modelling using shape priors. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, New York, NY*, 2006.
5. A. Del Bue, F. Smeraldi, and L. Agapito. Non-rigid structure from motion using ranklet–based tracking and non-linear optimization. *Image and Vision Computing*, 25(3):297–310, March 2007.
6. G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins, 1989.
7. R. I. Hartley and F. Schaffalitzky. Powerfactorization: an approach to affine reconstruction with missing and uncertain data. In *Australia-Japan Advanced Workshop on Computer Vision*, Adelaide, Australia, September 2003.
8. K. Kanatani and Y. Sugaya. Factorization without factorization: complete recipe. *Memories of the Faculty of Engineering, Okayama University*, 38(1–2):61–72, 2004.
9. Taeone Kim and Ki-Sang Hong. Estimating approximate average shape and motion of deforming objects with a monocular view. *International Journal of Pattern Recognition and Artificial Intelligence*, 19(4):585–601, 2005.
10. C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization approach. *International Journal of Computer Vision*, 9(2):137–154, 1992.
11. L. Torresani, D. Yang, E. Alexander, and C. Bregler. Tracking and modeling non-rigid objects with rank constraints. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, Kauai, Hawaii*, 2001.
12. J. Xiao, J. Chai, and T. Kanade. A closed-form solution to non-rigid shape and motion recovery. *International Journal of Computer Vision*, 67(2):233–246, April 2006.

| Missing % | Noise | | | | |
|-----------|-------|------|------|------|------|
|           | 0 | 0.5 | 1 | 1.5 | 2 |
| 10% | 1, 32 | 1, 47 | 1, 89 | 2, 11 | 2, 13 |
| 20% | 2, 85 | 3, 69 | 3, 45 | 3, 69 | 4, 05 |
| 30% | 3, 75 | 4, 74 | 4, 76 | 5, 03 | 5, 78 |
| 40% | 3, 99 | 4, 64 | 5, 18 | 5, 47 | 6, 87 |

Rotation Error

| Missing % | Noise | | | | |
|-----------|-------|------|------|------|------|
|           | 0 | 0.5 | 1 | 1.5 | 2 |
| 10% | 0.84 | 1.10 | 1.02 | 1.38 | 1.94 |
| 20% | 1.26 | 1.38 | 2.05 | 1.26 | 2.55 |
| 30% | 1.41 | 1.62 | 2.19 | 2.21 | 2.18 |
| 40% | 1.78 | 1.86 | 1.96 | 2.39 | 2.40 |

3-D Structure error

**Table 1.** Left: Mean of the the absolute rotation error expressed in degrees. Right: 3-D reconstruction error expressed in percentage relative to the scene size. The variance of the added noise is expressed in terms of image pixel. The value are computed on 10 trials for each configuration of noise and missing data ratios.
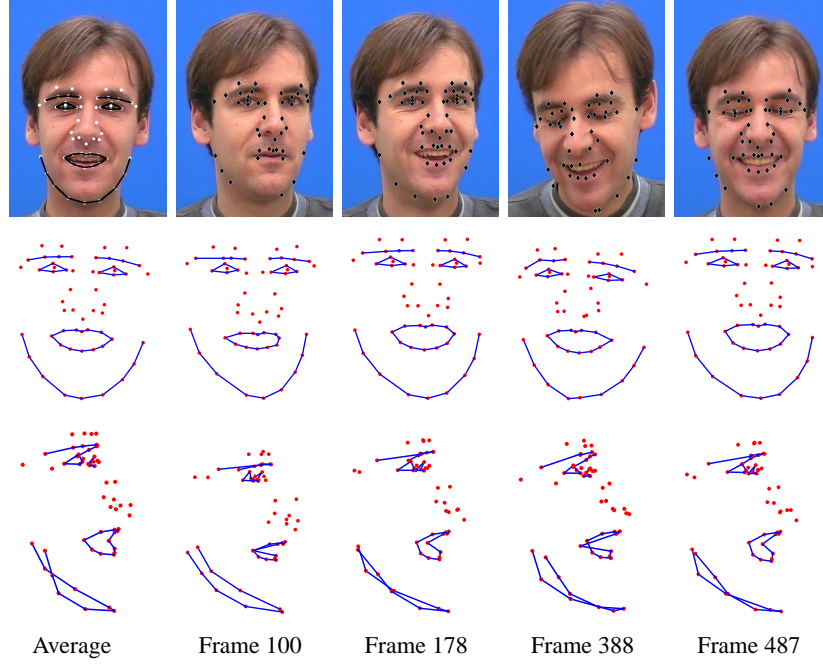


Average    Frame 100    Frame 178    Frame 388    Frame 487

**Fig. 1.** The first column shows the complete set of 56 points used for reconstruction (first row) and the recovered 3-D average shape (front and side views). The remaining columns present 4 key frames of the sequence with the available points at each frame. The second and third rows present respectively the front and side views of the reconstructed 3-D structure after non-linear optimisation. The number of basis shapes was fixed to $K = 6$.