

Benno Heigl, Heinrich Niemann:
Camera Calibration from Extended
Image Sequences for Lightfield
Reconstruction

Universität Erlangen–Nürnberg,
Lehrstuhl für Mustererkennung,
Martensstr. 3, 91058 Erlangen,
Germany

Email:

`heigl@informatik.uni-erlangen.de`

Workshop: Vision, Modeling and
Visualization
Erlangen, November 1999

Pages: 43–50

Abstract:

This contribution treats the problem of camera calibration of extended image sequences with the goal to generate lightfields. Our method is mainly based on the factorization method by Sturm and Triggs which provides a closed-form solution if all projections of all scene points into all cameras are known. We extend the original method for extended image sequences, especially when scene points appear and disappear in the projections of the sequence, which is natural when moving a camera around an object. We apply the original method on parts of the sequence and merge the results together. As our experiments show, in spite of these extension, our method produces comparable results to the original method.

Camera Calibration from Extended Image Sequences for Lightfield Reconstruction

Benno Heigl ^{*}, Heinrich Niemann

Universität Erlangen–Nürnberg, Lehrstuhl für Mustererkennung,
Martensstr. 3, 91058 Erlangen, Germany
Email: heigl@informatik.uni-erlangen.de

Abstract

This contribution treats the problem of camera calibration of extended image sequences with the goal to generate lightfields. Our method is mainly based on the factorization method by Sturm and Triggs which provides a closed-form solution if all projections of all scene points into all cameras are known. We extend the original method for extended image sequences, especially when scene points appear and disappear in the projections of the sequence, which is natural when moving a camera around an object. We apply the original method on parts of the sequence and merge the results together. As our experiments show, in spite of these extension, our method produces comparable results to the original method.

1 Introduction

The goal of our work is to construct a Lightfield [5] from an uncalibrated image stream which is taken with a hand-held camera. For a realistic visualization many views all around a scene have to be taken. Therefore, we need a stable method for calibrating especially long image sequences.

Many investigations have been made to solve the problem of camera calibration when just knowing projections of scene points with

unknown pose. Most approaches do a reconstruction from up to three successive images of the sequence and merge these reconstructions together by minimizing the total reprojection error of the reconstructed points (e.g. [1, 4]). They do not treat all data uniformly, but they sequentially reconstruct new camera poses with respect to those calibrated before. The whole reconstruction depends on the calibration of the first cameras.

For the case of orthographic projection, the first closed-form solution was given in [12]. There, a measurement matrix is built from projected image points and the camera positions as well as the scene points are reconstructed by a single factorization of this measurement matrix. Sturm and Triggs showed in [11] how to do a comparable solution for the case of perspective projection by building the measurement matrix from image projections and the corresponding projective depths. This method is very elegant as it recovers scene geometry and camera parameters for all images in one single step. The disadvantage of the closed-form solution is that it supposes that all scene points are visible throughout the whole image sequence. For real sequences this assumption does not hold.

In this contribution we focus on adapting the method [11] to sequences where points disappear and also new points appear. We apply the factorization method to parts of the sequence and merge them together. Doing this, we can benefit from the advantages of

^{*}This work is partially funded by the German Research Foundation (DFG) under grant number SFB 603

the closed-form solution and also we are able to treat the problem of loosing points.

In all these methods for recovering structure from motion, in a first step the reconstruction is performed up to an unknown projective transformation. By applying self-calibration methods, constraints on the cameras can be introduced to get a reconstruction up to an unknown similarity transformation. We apply the linear method [8] to do self-calibration, which is also needed to merge together the partial solutions.

This article is structured as follows. In section 2, we introduce some notations. Section 3 gives a short description of the original factorization method of Sturm and Triggs and section 4 shows briefly how to do self-calibration with the method [8]. In section 5 we show how to apply these methods to real sequences. As in this step not all projected points can be considered, section 6 shows how the additional information given by these points can be used. Section 7 gives an evaluation of our method by testing it on simulated data.

2 Notation

In the following we presume some knowledge of projective geometry. An excellent introduction can be found in [7] which is focused on the special requirements for image analysis.

All vectors are denoted bold face and describe homogeneous vectors if not stated otherwise. A 3-D scene point is denoted by 4-vector $\mathbf{p}_k = (p_{k1}, p_{k2}, p_{k3}, p_{k4})^T$ corresponding to the Euclidean 3-vector $(p_{k1}/p_{k4}, p_{k2}/p_{k4}, p_{k3}/p_{k4})^T$, if the point is finite. Similar, a 2-D image coordinate of the projection of \mathbf{p}_k into image number i is denoted by a homogeneous 3-vector $\mathbf{q}_{ik} = (q_{ik1}, q_{ik2}, q_{ik3})^T$. This implies that if a homogeneous vector is scaled by an arbitrary scalar not being zero, it does not change its Euclidean meaning.

The projection \mathbf{q}_{ik} of \mathbf{p}_k is achieved by a multiplication with a 3×4 projection matrix \mathbf{P}_i :

$$\lambda_{ik} \mathbf{q}_{ik} = \mathbf{P}_i \mathbf{p}_k \quad , \quad (1)$$

with $i = 1, \dots, m$ and $k = 1, \dots, n$. In a Euclidean framework, \mathbf{P}_i can be factorized as $\mathbf{P}_i = \mathbf{K}_i \mathbf{R}_i^T (\mathbf{I}_3 \mid -\mathbf{t}_i)$, with \mathbf{K}_i being an upper triangular calibration matrix and \mathbf{R}_i , \mathbf{t}_i being a rotation matrix and a (Euclidean!) translation vector transforming the world coordinate system into the appropriate camera coordinate system.

If a variable is marked with a prime (\mathbf{p}'), it describes an estimation for the true value. If a variable is marked with a hat ($\hat{\mathbf{p}}$), it describes a true value which is transformed projectively. If both marks are present ($\hat{\mathbf{p}}'$), a variable describes a projectively transformed estimation of the true value.

3 Factorization Method

In this section we give a brief description of the factorization method of Sturm and Triggs [11]. Their main idea is to build a *measurement matrix* \mathbf{W} and to solve the structure-from-motion problem by performing a single SVD (singular value decomposition, [9]) of \mathbf{W} .

To build the measurement matrix, the so-called *projective depths* λ_{ik} (see equation 1) have to be known. We suppose that all vectors \mathbf{q}_{ik} , \mathbf{p}_k are scaled such that the last component is 1 and all \mathbf{P}_k scaled such that their last row has norm 1. In this case the projective depths correspond to the orthogonal distances from the focal plane of each camera.

To recover these projective depths from projections only, fundamental matrices \mathbf{F}_{ij} and epipoles \mathbf{e}_{ij} between neighboring images i, j must be reconstructed. This can be achieved by applying known estimation techniques as described in [6]. In this contribution (following Sturm and Triggs) we apply the linear method of Hartley [3]. The projective depths can be updated from frame i to frame j by applying following equation:

$$\lambda_{ik} = \frac{(\mathbf{e}_{ij} \times \mathbf{q}_{ik})(\mathbf{F}_{ij} \mathbf{q}_{jk})}{\|\mathbf{e}_{ij} \times \mathbf{q}_{ik}\|^2} \lambda_{jk} \quad . \quad (2)$$

The projective depths of the first view can be chosen arbitrarily, e.g. $\lambda_{1k} = 1$.

Now all equations 1 for all scene points and all their projections are combined to one single matrix equation:

$$\mathbf{W} = \underbrace{\begin{pmatrix} \mathbf{P}_1 \\ \mathbf{P}_2 \\ \vdots \\ \mathbf{P}_m \end{pmatrix}}_{\Pi} \underbrace{(\mathbf{p}_1 \mathbf{p}_2 \dots \mathbf{p}_n)}_{\Psi} \quad \text{with} \quad (3)$$

$$\mathbf{W} = \begin{pmatrix} \lambda_{11} \mathbf{q}_{11} & \lambda_{12} \mathbf{q}_{12} & \dots & \lambda_{1n} \mathbf{q}_{1n} \\ \lambda_{21} \mathbf{q}_{21} & \lambda_{22} \mathbf{q}_{22} & \dots & \lambda_{2n} \mathbf{q}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{m1} \mathbf{q}_{m1} & \lambda_{m2} \mathbf{q}_{m2} & \dots & \lambda_{mn} \mathbf{q}_{mn} \end{pmatrix}.$$

The estimation \mathbf{W}' can be built knowing estimates for all projections of all scene points into all images and their corresponding projective depths. Note that the scale of each \mathbf{p}_k and the scale of each \mathbf{P}_i can be chosen arbitrarily if the λ_{ik} change accordingly. Therefore, to achieve a good numerical conditioning, each column c of $\mathbf{W}' = [w']_{rc}$ is rescaled such that $\sum_{r=1}^{3m} w'_{rc}^2 = 1$ and each triplet of rows $(3i - 2, 3i - 1, 3i)$ is scaled such that $\sum_{c=1}^n \sum_{l=3i-2}^{3i} w'_{cl}^2 = 1$.

An SVD is applied on \mathbf{W}' resulting in the factorization $\mathbf{W}' = \mathbf{U} \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_s) \mathbf{V}$. As the rank of \mathbf{W} in the noise-free case is 4, we can, for the case having noise, recover the best rank 4 approximation (in the sense of least squares) by composing $\mathbf{W}'' = \mathbf{U}' \text{diag}(\sigma_1, \sigma_2, \sigma_3, \sigma_4) \mathbf{V}'$ with \mathbf{U}' being the upper $3m \times 4$ sub-matrix of \mathbf{U} and \mathbf{V}' being the upper $4 \times n$ sub-matrix of \mathbf{V} . Knowing this factorization, estimates $\widehat{\mathbf{P}}'_i$ and $\widehat{\mathbf{p}}'_k$ can be found:

$$\widehat{\Pi}' = \mathbf{U}' \Sigma = \quad \text{and} \quad \widehat{\Psi}' = \Sigma \mathbf{V}' \quad (4)$$

with $\Sigma = \text{diag}(\sqrt{\sigma_1}, \sqrt{\sigma_2}, \sqrt{\sigma_3}, \sqrt{\sigma_4})$. The hats shall show that this representation is not unique as for an arbitrary 4×4 transformation matrix \mathbf{T} having rank 4 following equation holds: $(\widehat{\mathbf{P}}'_i \mathbf{T})(\mathbf{T}^{-1} \widehat{\mathbf{p}}'_k) = \widehat{\mathbf{P}}'_i \widehat{\mathbf{p}}''_k$, and similar $(\widehat{\Pi}' \mathbf{T})(\mathbf{T}^{-1} \widehat{\Psi}') = \widehat{\Pi}' \widehat{\Psi}''$. Therefore the structure of the scene and the projection matrices both can be reconstructed just up to an unknown projective transformation \mathbf{T} . This re-

striction is common to all reconstruction techniques if no additional knowledge on the properties of the cameras is available. To introduce such restrictions on the camera parameters, self-calibration methods can be applied as described in the following section.

4 Self-Calibration

As we have seen in the previous section, scene reconstruction can be done just up to an unknown projective transformation \mathbf{T} . The goal of self-calibration is to recover an estimation for matrix \mathbf{T}' which maps the reconstructed scene points to estimations of the real scene points: $\mathbf{p}'_k = \mathbf{T}'^{-1} \widehat{\mathbf{p}}'_k \approx \mathbf{p}_k$ and accordingly $\mathbf{P}'_i = \widehat{\mathbf{P}}'_i \mathbf{T}' \approx \mathbf{P}_i$. This can be performed only by imposing some restrictions on the camera parameters.

In [14] the concept of the *absolute quadric* is introduced for doing self-calibration. The absolute quadric is described by the 4×4 matrix $\Omega = \begin{pmatrix} \mathbf{I}_{3 \times 3} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}$. We restrict to an algebraic description and omit here the geometric interpretation of this concept as it is insubstantial for performing self-calibration. It is obvious that following equation holds:

$$\mathbf{P}_i \Omega \mathbf{P}_i^T = \rho \mathbf{K}_i \mathbf{K}_i^T, \quad (5)$$

where \mathbf{K}_i denotes the upper triangular calibration matrix. Suppose that the true scene points and camera matrices are transformed by \mathbf{T} . Then we see from Eq. 5 that Ω is transformed to $\widehat{\Omega} = \mathbf{T} \Omega \mathbf{T}^T$. If \mathbf{T} is a similarity transformation (meaning that $\mathbf{T} = \begin{pmatrix} \sigma \mathbf{R} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{pmatrix}$ with \mathbf{R} being a rotation matrix), the absolute quadric is invariant, therefore $\widehat{\Omega} = \Omega$. But in our case, \mathbf{T} is an arbitrary projective transformation, and in general $\widehat{\Omega} \neq \Omega$. In [8] a linear method is presented to recover an estimate $\widehat{\Omega}'$ uniquely for a given set of projection matrices, if all calibration matrices

have the form $\mathbf{K}_i = \begin{pmatrix} f_i & 0 & 0 \\ 0 & f_i & 0 \\ 0 & 0 & 1 \end{pmatrix}$. When

the image skew, the principal point, and the aspect ratio is known, this can be achieved by transforming the measured image coordinates so that the principal point is $(0,0)$, the aspect ratio is 1, and the skew is 0. The only unknown camera parameter then is the focal length. The idea is to transform $\widehat{\mathbf{P}}'_1$ to be $(\mathbf{I}_{3 \times 3} \mathbf{0})$:

$$(\mathbf{I}_{3 \times 3} \mathbf{0}) = \widehat{\mathbf{P}}'_1 \underbrace{\begin{pmatrix} \widehat{\mathbf{P}}'_1 \\ \mathbf{0}^T \mathbf{1} \end{pmatrix}^{-1}}_{\mathbf{T}_n}. \quad (6)$$

All other projection matrices are transformed accordingly by multiplying $\widehat{\mathbf{\Pi}}'$ with \mathbf{T}_n . (For consistency, $\widehat{\mathbf{\Psi}}'$ is pre-multiplied with \mathbf{T}_n^{-1} .) In the noise-free case, following equation holds:

$$\mathbf{A}_i = \rho \begin{pmatrix} f_i^2 & 0 & 0 \\ 0 & f_i^2 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \widehat{\mathbf{P}}_i \widehat{\mathbf{\Omega}} \widehat{\mathbf{P}}_i^T. \quad (7)$$

We see that these four equations have to be satisfied: $A_{i11} = A_{i22}$, $A_{i12} = A_{i13} = A_{i23} = 0$. Therefore having m frames, we get $4(m-1)$ equations for estimating the symmetric absolute quadric $\widehat{\mathbf{\Omega}}'$, which is parameterized by the 10 elements of its upper triangular sub-matrix. The focal length of each camera then can be calculated from \mathbf{A}_i . From $\widehat{\mathbf{\Omega}}'$, the plane at infinity can be retrieved as the nullspace of $\widehat{\mathbf{\Omega}}'$ (e.g. using SVD), and is denoted by the row-vector \mathbf{p}'_∞ . Supposing the camera coordinate system to coincide with the world coordinate system, we know that the unknown matrix \mathbf{T}' must transform $(\mathbf{I}_{3 \times 3} \mathbf{0})$ to $(\mathbf{K}'_1 \mathbf{0})$. We also know that \mathbf{p}'_∞ must be transformed to the standard plane at infinity (0001) . Applying these restrictions, the matrix \mathbf{T}' is determined uniquely:

$$\mathbf{T}' = \begin{pmatrix} \mathbf{K}'_1 & \mathbf{0} \\ -\bar{\mathbf{p}}'_\infty \cdot \mathbf{K}'_1 & 1 \end{pmatrix}, \quad (8)$$

with $\bar{\mathbf{p}}'_\infty$ being the 3-vector consisting of the first three components of \mathbf{p}'_∞ divided by its fourth component. The scale of the last row of \mathbf{T} can be chosen freely. This is equal to scaling the whole scene with a unique factor which remains unknown.

5 Application to Real Image Sequences

The main disadvantage of the factorization method (section 3) is that all scene points must be visible throughout the whole sequence. In the analysis of real image sequences, this restriction is not practicable, because scene parts appear and disappear in dependence on the actual viewpoint and also there is no tracking algorithm which guarantees not to lose any scene points although they are visible. Hence we have to modify the algorithm to apply it to real environments.

Factorization of windows. During tracking point features in extended image sequences (for example, using the method [10]), it is usual to select a given number of point features (e.g. 1000) and to track them as long as possible. If in a particular frame features are lost, new features are selected to complete the number of tracked features. Therefore the factorization method cannot be applied to the whole sequence. We can inspect “windows” of frames where enough features could be tracked completely. On the one hand, the size of each window should not be too small, as the reconstruction error reduces with increasing number of frames (compare to the results in [11]). On the other hand, the window should not be too large, as the number of points which could be tracked throughout the whole window decreases and therefore the reconstruction is more sensitive to errors in the trajectory corresponding to a single scene point. The windows can be chosen to be a complete partitioning of the frames or they can be chosen to overlap. In the latter case, also the camera centers of the overlapping frames can be matched to each other. This is reasonable especially then, when the scene is compact and the cameras move around the scene with a large distance. Which choice performs best, depends on the particular situation. Let $\widehat{\mathbf{\Psi}}_v$ denote the matrix of the transformed scene points for window number v , according to the notation used in equation 3. \mathcal{I}_v denotes the set of corresponding indices

of scene points which are used for this reconstruction in window.

Merging windows. If we have calculated a projective reconstruction for each window, we must link these reconstructions to each other. As each reconstruction is done up to an unknown projective transformation, there also exists an unknown transformation \mathbf{T}_{vw} between two windows v and w which maps the transformed scene points of window v to the corresponding scene points of window w . Therefore these scene points must be visible in both windows. Let $\mathcal{I}_{vw} = \mathcal{I}_v \cap \mathcal{I}_w$ denote the set of indices of scene points which are visible in both windows v and w . Formally, for each scene point with an index $k \in \mathcal{I}_{vw}$, the matrix \mathbf{T}_{vw} should satisfy the equation

$$\rho {}^w\hat{\mathbf{p}}_k = \mathbf{T}_{vw} {}^v\hat{\mathbf{p}}_k, \quad (9)$$

where ${}^v\hat{\mathbf{p}}_k$ denotes the scene point k in the v -th window. It is a column vector of $\widehat{\Psi}_v$. If the vectors ${}^w\hat{\mathbf{p}}_k$ and ${}^v\hat{\mathbf{p}}_k$ are normalized such that their last component is 1, the scale ρ and the fourth component of the vector equation 9 can be eliminated. Therefore each correspondence of estimated scene points leads to three rows of a matrix:

$$\begin{pmatrix} {}^w\hat{\mathbf{p}}_k'^T & \mathbf{0}^T & \mathbf{0}^T & -{}^v\hat{\mathbf{p}}_{k1}' & {}^w\hat{\mathbf{p}}_k'^T \\ \mathbf{0}^T & {}^w\hat{\mathbf{p}}_k'^T & \mathbf{0}^T & -{}^v\hat{\mathbf{p}}_{k2}' & {}^w\hat{\mathbf{p}}_k'^T \\ \mathbf{0}^T & \mathbf{0}^T & {}^w\hat{\mathbf{p}}_k'^T & -{}^v\hat{\mathbf{p}}_{k3}' & {}^w\hat{\mathbf{p}}_k'^T \end{pmatrix}.$$

The approximative nullspace of this matrix can be solved using SVD and describes the vector being the concatenation of the column vectors of \mathbf{T}'_{vw} .

In the noise-free case, this merging procedure works. But having errors in the reconstruction, the result is numerically not stable, because there are minimized distances within a projectively transformed scene. As tiny distances may become huge after a projective transformation and vice versa, the particular distances don't have a meaningful sense. Hence we have to minimize these distances within an Euclidean framework; we have to apply the self-calibration before merging.

In general, it is enough to do the self-calibration for reconstruction of the first window and to calculate the transformation from

each window to the previous one. With this procedure the Euclidean framework is updated from one window to the next. But if there are small errors in the self-calibration of the first frame, the whole reconstruction will depend on this error. A final self-calibration step may be applied to improve the results.

Post-processing. The approach described above retrieves each camera matrix \mathbf{P}' up to an unknown scale factor ρ . As \mathbf{P} can be represented by \mathbf{K} , \mathbf{R} , and \mathbf{t} having fixed scales, we have to determine ρ to convert \mathbf{P}' to this representation: $\mathbf{P}' \rightarrow \mathbf{P}'/\rho$. For each projection matrix \mathbf{P}' , we can determine the oriented optical axis $\mathbf{a}' = (P'_{31} \ P'_{32} \ P'_{33})$ which is a Euclidean directional vector and is collinear with the third row of matrix \mathbf{R} . Therefore we know that $\|\rho\| = \|\mathbf{a}'\|$. The optical center \mathbf{t}' can be determined uniquely from \mathbf{P}' . As the reconstructed points must lie in front of the camera, the scalar product of the directional vector \mathbf{a}' with the vector $\mathbf{p}'_k - \mathbf{t}'$ should be positive for any visible scene point number k . To achieve this, we can determine $\rho = \text{sgn}(\mathbf{a}'^T(\mathbf{p}'_k - \mathbf{t}'))\|\mathbf{a}'\|$.

6 Considering All Points

Up to now we have reconstructed all projection matrices and the 3-D position of some scene points. But there remain points which are visible in fewer frames than the window lengths and therefore have not been considered yet. Knowing all projection matrices, we are able to reconstruct these scene points. We can again formulate this request as a linear optimization problem. Knowing the projection $\mathbf{q}_{ik} = (q_{ik1}, q_{ik2}, q_{ik3})^T$ of the scene point number k in the image i , following equation holds:

$$\left(\begin{pmatrix} q_{ik1} \\ q_{ik2} \end{pmatrix} \boldsymbol{\psi}_{i3} - \begin{pmatrix} \psi_{i1} \\ \psi_{i2} \end{pmatrix} \right) \mathbf{p}_k = 0, \quad (10)$$

where $\boldsymbol{\psi}_{ij}$ denotes the j -th row-vector of matrix \mathbf{P}_i . The scale of \mathbf{q}_{ik} has to be chosen such that its third component is 1. Writing this equation for each estimated projection of an unknown scene point and concatenating the

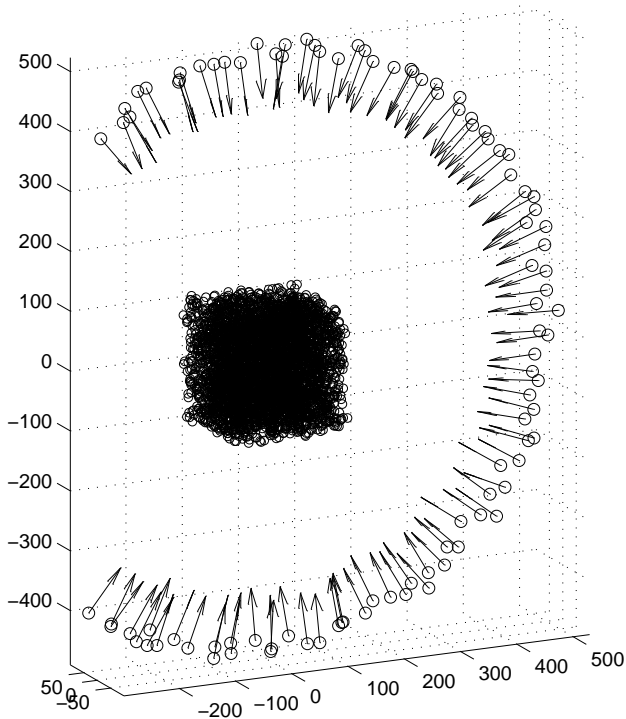


Figure 1: Example for a simulated test configuration. The arrows show the camera centers with their viewing directions; the points in the middle show the scene points. Note the small displacements of neighboring cameras.

left matrix of upper equation, we get a matrix whose nullspace describes the coordinates of the searched scene point \mathbf{p}'_k . The approximative nullspace again can be determined by applying SVD.

Another way to consider the remaining points is to complete the measurement matrix during each factorization step. As we have estimated the fundamental matrices between succeeding frames, we can predict their projections from three succeeding views using the method [2]. For this approach, many degenerate and nearly degenerate cases exist. E.g. the projection centers of three views must not be collinear, or the corresponding scene point must not lie on the plane built by the three projection centers. Therefore, the better way is to estimate the trifocal tensor (e.g. by method [13]) and to use that for prediction.

As small errors in the measured image coordinates as well as in the estimated funda-

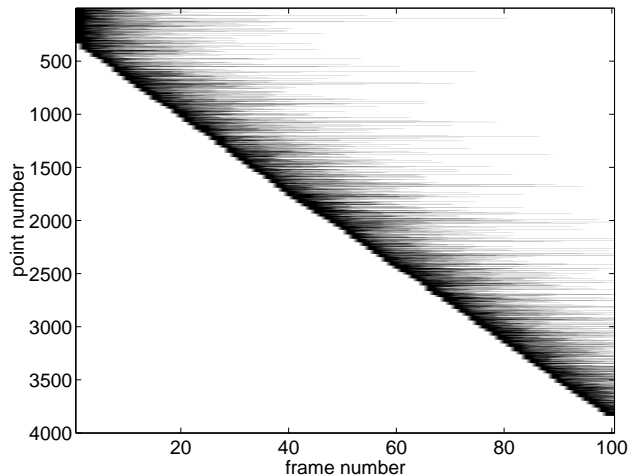


Figure 2: Simulated visibilities

mental matrices or trifocal tensors affect the predictions, this method of completing the measurement matrix should just be applied to fill small gaps. Otherwise, the reconstruction will depend too much on these erroneous estimated predictions.

7 Experiments

Simulation of camera movement and scene points. The scene points have been chosen to be randomly distributed in a cube with an edge length of 200. The camera moves around the center of the cube with a radius of 500. Its position as well as its orientation is perturbed randomly to avoid degeneracies of self-calibration (see section 4). This perturbation is also natural for sequences taken with a hand-held camera. The focal lengths have been chosen randomly in the range 800...1200. Figure 1 shows an example configuration. The scene points have been projected into each camera. To each projection, an error has been added with a uniform distribution within the range $-r \dots r$. We call r the *noise level*.

Simulation of visibility. To simulate the visibility of each point, we made the assumption that the probability p for losing a point in a frame is constant. Then we get a geometric distribution for the length t of the visibility for each point. Following equation

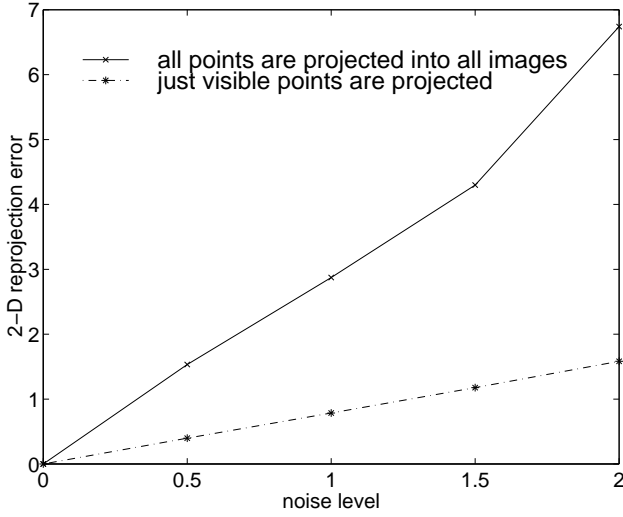


Figure 3: 2-D mean absolute pixel error of reprojections. These results were obtained for a window size of 20 frames with 5 overlapping frames between neighboring windows.

is used to create this distribution out from a uniform distribution within the range $0 \dots 1$: $t(x) = \log(1 - x)/\log(1 - p)$. We also made the assumption that a fixed number of points must be seen in each frame (in our example 400). This means that for each lost point a new one is introduced. Figure 2 shows an example of the simulated visibility in the case of $p = 0.1$.

Reprojection error. One measure for the evaluation of errors is the distance between the projections of the reconstructed points into the reconstructed cameras and the original projections. Figure 3 shows the mean absolute error in the reprojections measured in pixels. It can be seen that for each reconstructed camera, the reprojected image coordinates of those points which are visible in this camera have an error in the magnitude of the noise (dotted curve). If even those points are projected into a particular camera which are just visible in other cameras, the error increases (solid curve). This increasing error is caused by errors during merging together the reconstructions of neighboring windows.

3-D error of the reconstruction. Another accuracy measure is the 3-D error of the reconstructed points. Since the reconstruction is done just up to an unknown similar-

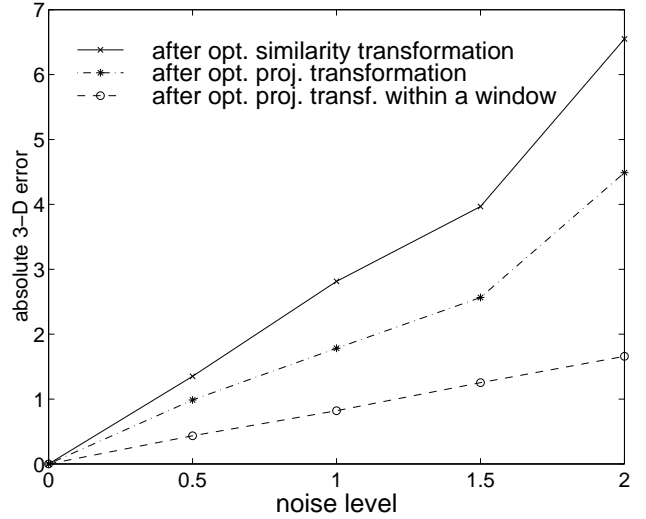


Figure 4: Mean absolute 3-D error for the same configuration in figure 3. It describes the mean difference between the true scene points and their corresponding reconstructed points after a specified transformation.

ity transformation, we have to determine that transformation which minimizes the 3-D error between the transformed reconstructed points and the true scene points.

Figure 4 shows the absolute 3-D errors in units of the simulated 3-D structure after that optimal similarity transformation (solid curve). The errors after transforming with an optimal projective transformation are less. This indicates a projective skew in the reconstruction. By applying a non-linear refinement of the self-calibration step, this skew should be reduced. If the optimal projective error is determined separately within each window, the 3-D error is even less. As in the 2-D case, this again indicates errors during the merging step.

Choice of the number of overlapping frames. As described in section 5, neighboring windows may be chosen to overlap by some frames to increase stability. Figure 5 shows the 3-D error after the optimal similarity transformation in dependence on the number of overlapping frames. It can be seen that a significant improvement can be noticed just for higher noise levels.

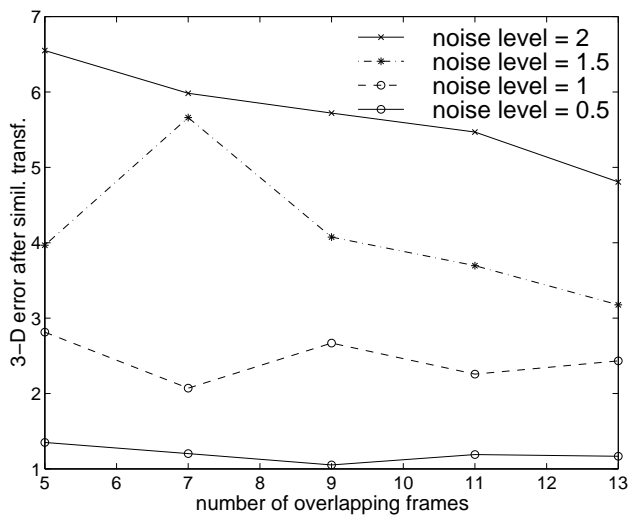


Figure 5: Error in dependence on the number of overlapping frames for different noise levels.

8 Conclusion

In this contribution we have shown how to apply the factorization method of Sturm and Triggs to real image sequences. We have seen that it is necessary to merge together partial solutions within a Euclidean framework which can be recovered by applying a self-calibration method. The experiments show that our method performs well, but improvements could be achieved by enhancing the merging step.

Acknowledgements. We would like to thank Marc Pollefeys from ESAT-PSI of the K. U. Leuven for giving us hints how to perform the self-calibration step.

References

- [1] P. A. Beardsley, P. H. S. Torr, and A. Zisserman. 3D model acquisition from extended image sequences. In *Proceedings ECCV*, pages 683–695, 1996.
- [2] O. Faugeras and L. Robert. What can two images tell us about a third one? Technical report, INRIA, Sophia Antipolis, 1993. technical report.
- [3] R. I. Hartley. In defense of the eight-point algorithm. *PAMI*, 19(6):580–593, 1997.
- [4] R. Koch, M. Pollefeys, and Luc Van Gool. Multi viewpoint stereo from uncalibrated video sequences. In *Proceedings ECCV*, pages 55–71. Springer, 1998.
- [5] M. Levoy and P. Hanrahan. Light field rendering. In *Proceedings SIGGRAPH*, pages 31–45, 1996.
- [6] Q.-T. Luong, R. Deriche, and O. Faugeras. On determining the fundamental matrix: Analysis of different methods and experimental results. Technical report, INRIA, Sophia Antipolis, 1993. technical report.
- [7] R. Mohr and B. Triggs. Projective geometry for image analysis. In *Int. Symp. Photogrammetry and Remote Sensing*, July 1996. Tutorial.
- [8] M. Pollefeys, R. Koch, and L. Van Gool. Self-calibration and metric reconstruction in spite of varying and unknown internal camera parameters. In *Proceedings ICCV*, pages 90–95, Bombay, 1998.
- [9] W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling. *Numerical Recipes in C – The Art of Scientific Computing*. Cambridge University Press, New York, 1990.
- [10] J. Shi and C. Tomasi. Good features to track. In *Proceedings CVPR*, pages 593–600. IEEE Computer Society Press, June 1994.
- [11] P. Sturm and B. Triggs. A factorization based algorithm for multi-image projective structure from motion. In *Proceedings ECCV*, pages 709–720. Springer, 1996.
- [12] C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: a factorization method. *International Journal of Computer Vision*, 9(2):137–153, 1992.
- [13] P. H. S. Torr and A. Zisserman. Robust parameterization and computation of the trifocal tensor. *Image and Vision Computing*, 15:591–605, 1997.
- [14] B. Triggs. Autocalibration and the absolute quadric. In *Proceedings CVPR*, pages 609–614, June 1997.