



HAL
open science

Automatic camera calibration using multiple sets of pairwise correspondences

Francisco Vasconcelos, Joao Pedro Barreto, Edmond Boyer

► **To cite this version:**

Francisco Vasconcelos, Joao Pedro Barreto, Edmond Boyer. Automatic camera calibration using multiple sets of pairwise correspondences. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018, 40 (4), pp. 791-803. 10.1109/TPAMI.2017.2699648 . hal-01675686

HAL Id: hal-01675686

<https://inria.hal.science/hal-01675686>

Submitted on 5 Jan 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Automatic camera calibration using multiple sets of pairwise correspondences

Francisco Vasconcelos, João P. Barreto, and Edmond Boyer

Abstract—We propose a new method to add an uncalibrated node into a network of calibrated cameras using only pairwise point correspondences. While previous methods perform this task using triple correspondences, these are often difficult to establish when there is limited overlap between different views. In such challenging cases we must rely on pairwise correspondences and our solution becomes more advantageous. Our method includes an 11-point minimal solution for the intrinsic and extrinsic calibration of a camera from pairwise correspondences with other two calibrated cameras, and a new inlier selection framework that extends traditional RANSAC to sampling across multiple datasets. Our method is validated on different application scenarios where a lack of triple correspondences might occur: addition of a new node to a camera network; calibration and motion estimation of a moving camera inside a camera network; and addition of views with limited overlap to a Structure-from-Motion model.

Index Terms—Camera Calibration, Camera Networks, Minimal Algorithms, RANSAC



1 INTRODUCTION

A camera network, in the context of this article, is a set of cameras with synchronous image acquisition and partial overlap in the field-of-views (FOVs). These camera networks are popular in application domains that are concerned with the capture, the record, and the analysis of dynamic scenes, such as surveillance, gait analysis, human-motion capture, or 3D modelling for the movie industry [1]. Such applications invariably require the camera network to be calibrated, meaning that both intrinsic and extrinsic parameters must be known for all camera nodes in order to fuse the multiple-view information.

The problem of camera network calibration has been broadly addressed in the past. One possibility is to use a known calibration object, such as a checkerboard pattern, that is simultaneously observed by all nodes [2], [3], [4]. Some authors have recently proposed to observe the object through planar mirror reflections in order to handle situations of little or no overlap in the FOVs [5], [6]. Another option is to freely move a Light-Emitting Diodes (LED) in a dark room for obtaining accurate image correspondences that are used as inputs into factorization step [7], [8], [9]. All these calibration procedures are explicit, in the sense that they require substantial human intervention, and are meant to be carried as an initial off-line step before

starting operating the network.

In spite of the many explicit methods for accomplishing camera network calibration, there are situations for which an automatic, unsupervised scheme is highly advantageous such as adding or adjusting the position of a camera while the network is in operation. Some efforts have been made to accomplish this task in real time in the context of sports broadcasting [10], [11]. These approaches take advantage from the fact that a sports field provides easily detectable features from a planar region. When no assumptions are made about the viewed scene a more suitable alternative is to estimate the camera parameters from natural image point correspondences with neighbouring calibrated views. It is well known that the camera projection matrix can be estimated in a DLT like manner from 6 or more triple correspondences [12], [13], with triple correspondence standing for an image point that is viewed in two other calibrated nodes such that it is possible to find its 3D coordinates. Unfortunately, and as shown in Fig. 1, triple correspondences are often difficult to establish in practice, either because cameras are separated by a wide-baseline and present very different perspectives, or because the dynamic scene creates relative occlusions that preclude matching. Thus, we propose to relax the requirements in the input data and perform the calibration from independent pairwise correspondences. The problem has 11 unknown parameters (5 intrinsics and 6 extrinsics), which means that in theory the solution can be fully constrained from a total of 11 pairwise correspondences between the uncalibrated camera and two distinct calibrated views. We propose a minimal solution for the problem (the 11-pt algorithm) that requires 7 matches with the first view and 4 matches with the second view.

- *F. Vasconcelos is with University College London, United Kingdom. This work was carried while he was with Institute for Systems and Robotics, University of Coimbra, Portugal*
- *J. P. Barreto is with the Institute for Systems and Robotics, University of Coimbra, Portugal*
- *E. Boyer is with INRIA Grenoble Rhône-Alpes, France*

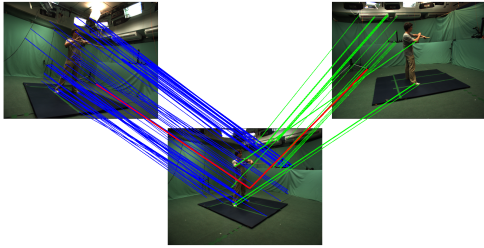


Fig. 1. Correspondences extracted from SIFT features. Given the wide baseline between the views there is a single reliable triple correspondence (red) while there are many reliable pairwise correspondences (blue and green).

The present article builds on our previous conference publication [14] that discloses the 11-pt algorithm. It is important to note that [14] does not include any experimental validation on datasets contaminated by outliers, since using standard RANSAC methods [15] proved to be unsuccessful. The standard RANSAC formulation assumes a single set of correspondences that is iteratively sampled to compute candidate models that are used to split data into inliers and outliers. In our case the search requires sampling more than one dataset that might present different inlier-outlier statistics. It is shown that overlooking this fact and applying the standard RANSAC formulation leads to poor results. Thus, we propose a modified RANSAC version specifically designed to simultaneously sample multiple datasets. The usefulness of such RANSAC extension goes beyond the calibration problem at hands and can benefit other algorithms such as Structure-from-Motion (SfM) using stereo cameras [16], [17], multi-view camera rigs [18], [19], or a mixture of point and plane correspondences [20], [21]. The contributions in this paper can be summarised as follows:

- A minimal algorithm for estimating the intrinsic and extrinsic parameters of a camera from 11 independent pairwise correspondences with two other calibrated cameras.
- Extensions of the well known RANSAC [15], MLESAC [22], and MAPSAC [23] formulations for sampling multiple different datasets.
- A simple and efficient implementation of the complete calibration solution that is able to calibrate a camera from correspondences with two or more cameras. This implementation is tested in three different scenarios: calibrating stationary camera nodes in a network; finding the parameters of a hand-held camera that freely moves within a camera network space to acquire close-ups of foreground dynamic scenes; Adding new viewpoints to a SfM model. The experiments demonstrate that in challenging scenarios with limited overlap in fields of view or significant

change in viewpoint our method is able to outperform the state-of-the-art.

2 RELATED WORK

Since a camera array or network can be understood as a generalised camera [18], the extrinsic calibration of a camera from independent pairwise correspondences with multiple views relates with the problem of relative pose estimation between non-central cameras. It is well known that the rotation and translation between two generalised cameras can be solved linearly from 17 correspondences [24] and solved minimally from 6 pairwise correspondences [19]. However, these methods degenerate in many particular configurations, namely when one of the generalised views is a pin-hole as it happens in our case. In the case of the camera network having just two nodes the extrinsic calibration problem from pairwise correspondences can be potentially solved using methods developed for visual odometry using stereo cameras. There is a minimal solution for the relative pose between stereo pairs using 6 pairwise correspondences [16] that estimates an up-to scale relative pose solution using 5 correspondences with one camera [25] and solves the scale factor with an additional correspondence from another camera. A non-minimal solution using 10 correspondences was also proposed for the case of any arbitrary combination of correspondences between the 4 views of two stereo rigs [17]. Since in this paper we focus on both intrinsic and extrinsic calibration from pairwise correspondences, the above mentioned works relate but do not directly apply.

Whenever triple correspondences are available the calibration objective can be accomplished using standard techniques described in text books [12], [13]. These approaches typically rely on reconstructing 3D points from the the calibrated stereo views via triangulation [26], and using these points as known reference to calibrate the third view [8], [13]. Unfortunately, and as discussed in the introduction, triple correspondences are not always available. A possible alternative is to build a measurement matrix with the image correspondences, and perform projective factorization using the Sturm-Triggs algorithm [27] with a suitable extension for handling missing data [7]. However, this class of methods is meant for problems with multiple cameras and large number of correspondences, and it is unlikely that the approach will converge to a solution using only pairwise correspondences. Levi and Werman propose to build a viewing graph where pairs of camera nodes are linked by their fundamental matrices [28]. They show that, given a subset of known fundamental matrices, it is possible to determine the remaining edges in the graph as far as each camera node is connected with at least two other camera nodes. This condition is not verified whenever we aim to calibrate a camera that has been just added

to an existing camera network. In [29] Josephson et al. investigate the problem of calibrating a camera node from mixtures of triple and pairwise correspondences. However, and to the best of our knowledge, the calibration of a camera using exclusively independent pairwise correspondences with two other views has never been addressed in the literature before.

Additionally we are interested in robust estimation with RANSAC when candidate solutions are generated by sampling multiple datasets. This problem is only briefly addressed in [16] when estimating the relative pose between stereo rigs. In that problem RANSAC must independently select one sample from 3 datasets containing 2-view, 3-view, and 4-view correspondences. It is shown that the number of RANSAC iterations must be computed in a different way when the different datasets have different inlier ratios, and also that different types of correspondences should be weighted differently on the cost function. However, the observations made in [16] can only be directly extended to problems where we know exactly how many samples are selected in each dataset. In our calibration problem this might not be the case, since there are different combinations of pairwise correspondences among different cameras that can generate a solution. Adapting other variants from the RANSAC-family (e. g. MLESAC [22]) to a multiple dataset framework have also never been addressed before.

2.1 Notation

Scalars are represented by plain letters, e.g. λ , vectors are indicated by bold symbols, e.g. \mathbf{t} , and matrices are denoted by letters in sans serif font, e.g. \mathbf{T} . 3D lines are expressed in homogeneous Plucker coordinates, e.g. the 6×1 vector \mathbf{L} . The equality up to scale is denoted by \sim in order to be distinguished from the strict equality $=$, and the operator $[\mathbf{v}]_{\times}$ designates the 3×3 skew symmetric matrix of a 3×1 vector \mathbf{v} . We also use matrix superscripts, e. g. $\mathbf{T}^{\{n\}}$, to denote its n th column. \otimes denotes the kronecker product.

3 PROBLEM STATEMENT

Let us consider two calibrated cameras C_A and C_B , such that the matrices of intrinsic parameters are K_A and K_B , and the absolute poses are expressed in a world coordinate system O_w by the rotation matrices R_A and R_B , and the translation vectors \mathbf{t}_A and \mathbf{t}_B . Consider an additional camera C for which both the intrinsic calibration K , and the extrinsic calibration R , \mathbf{t} are unknown. Our article addresses the problem of calibrating this third camera using as input data a set of a image correspondences $(\mathbf{x}^{(i)}, \mathbf{x}_A^{(i)})$ between C and C_A , and set of b image correspondences $(\mathbf{x}^{(a+j)}, \mathbf{x}_B^{(j)})$ between C and C_B (Fig. 2). We assume that the two

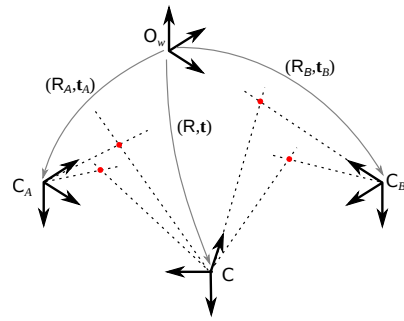


Fig. 2. We consider the problem of fully calibrating the camera C , given pairwise correspondences with two calibrated cameras C_A and C_B .

sets of pairwise matches are independent, meaning that

$$\mathbf{x}^i \neq \mathbf{x}^{a+j}, \forall i=1\dots a, j=1\dots b. \quad (1)$$

4 LINEAR CONSTRAINTS

In this section we derive a system of linear equations that has a minimum number of unknowns and fully constrains the camera calibration. The problem is formulated in the context of epipolar geometry between general camera models [18], with one side being the uncalibrated pin-hole camera C , and the other side being the pair of calibrated cameras C_A and C_B that can be understood as a particular instance of a non-central imaging device denoted by $C_A \cup C_B$. It is shown below that under such configuration the corresponding back-projection lines must satisfy a bilinear relation expressed by a 3×5 matrix, and that the estimation of the epipolar geometry using a DLT-like approach cannot be achieved with less than 14 pairwise matches.

Note that when the intrinsics are known, this problem is a particular case of the pose estimation between calibrated general camera models [18] that has already been solved both linearly [24] and using the minimal number of 6 pairwise correspondences [19].

4.1 Line Incidence Relations

Let \mathbf{x}_A and \mathbf{x}_B be image points in C_A and C_B . Since the cameras are fully calibrated, the corresponding back-projection lines \mathbf{L}_A and \mathbf{L}_B can be expressed in the common world reference frame O_w by a homogeneous Plucker vector [30]

$$\mathbf{L}_{A/B} \sim \begin{pmatrix} \mathbf{d}_{A/B} \\ \mathbf{m}_{A/B} \end{pmatrix}, \quad (2)$$

with the 3-vectors $\mathbf{d}_{A/B}$ and $\mathbf{m}_{A/B}$ being respectively the direction and the moment of the line. In a similar manner, an image point \mathbf{x} in C gives rise to a back-projection line \mathbf{L} that is represented in the local camera reference frame by

$$\mathbf{L} \sim \begin{pmatrix} \mathbf{d} \\ \mathbf{0} \end{pmatrix}, \quad (3)$$

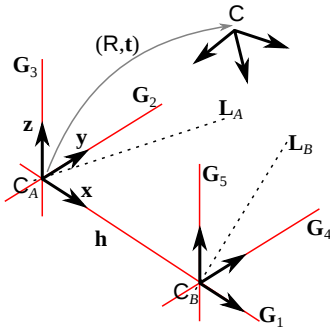


Fig. 3. The space generated by two bundles of lines (the rays of 2 pinhole cameras) can be fully represented as the linear span of $\{G_1, G_2, G_3, G_4, G_5\}$.

with the direction depending on the matrix of intrinsic parameters K

$$\mathbf{d} \sim K^{-1} \mathbf{x}. \quad (4)$$

If \mathbf{x} and $\mathbf{x}_{A/B}$ are image correspondences, then the back-projection lines \mathbf{L} and $\mathbf{L}_{A/B}$ must be incident. Given the rigid displacement between the reference frames O_w and C , and the condition for two lines in Plücker coordinates to intersect, then the following condition must hold [30]

$$\mathbf{L}^T \begin{pmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{R} & 0 \\ [\mathbf{t}]_{\times} \mathbf{R} & \mathbf{R} \end{pmatrix} \mathbf{L}_{A/B} = 0. \quad (5)$$

Since the moment of \mathbf{L} is always zero, then the above equation can be re-written as

$$\mathbf{d}^T ([\mathbf{t}]_{\times} \mathbf{R} \ \mathbf{R}) \mathbf{L}_{A/B} = 0. \quad (6)$$

Equation 6 is the particular case of the generalized epipolar constraint proposed in [18] when one of the cameras is a conventional pin-hole. However, and similarly to the general case, the bilinear relation between back-projection lines is expressed by a 3×6 matrix that encodes the calibration parameters. In a first glance it might seem that linearly estimating the 18 entries of this matrix up to a global scale factor can be carried with 17 or more image correspondences between C and the camera pair $C_A \cup C_B$. However, and as discussed below, these correspondences only provide 15 independent linear constraints in the matrix parameters.

In our case the parametrization of equation 6 leads to a linear estimation problem that is underdetermined. This is a situation similar to the degenerate configurations recently reported in [24] in the context of motion estimation using a calibrated multi-camera rig.

4.2 Compact linear formulation

The image rays belonging to two pinhole cameras C_A and C_B define a subset of lines that intersect a common axis \mathbf{h} (Fig 3). This subset is called a *linear line congruent* [31], and all its elements can be defined

as a linear combination of five lines G_1, G_2, G_3, G_4, G_5 that intersect \mathbf{h} . In our calibration problem every possible back-projection line $\mathbf{L}_{A/B}$ must intersect the line going through C_A and C_B (the baseline). Thus, the lines $\mathbf{L}_{A/B}$ can be represented in a unique manner as the linear combination of any 5 lines G_i that intersect the baseline

$$\mathbf{L}_{A/B} \sim \underbrace{\begin{pmatrix} G_1 & G_2 & G_3 & G_4 & G_5 \end{pmatrix}}_G \lambda_{A/B}, \quad (7)$$

where G is a 6×5 matrix with full rank, and $\lambda_{A/B}$ is a 5-vector defined up to scale. Replacing in equation 6 yields

$$\mathbf{d}^T ([\mathbf{t}]_{\times} \mathbf{R} \ \mathbf{R}) G \lambda_{A/B} = 0. \quad (8)$$

We have just re-written the epipolar constraint of equation 6 as a bilinear relation between the direction \mathbf{d} of the line \mathbf{L} in camera C , and the representation $\lambda_{A/B}$ of the back-projection line $\mathbf{L}_{A/B}$ in the generalized camera $C_A \cup C_B$. Since the bilinear relation is now encoded by a 3×5 matrix with 15 entries, then 14 image point correspondences are sufficient for estimating the epipolar geometry in a DLT-like manner.

Given the two arbitrary calibrated cameras, it is always possible to perform a change of reference frames for achieving the configuration exhibited in Fig. 3. We consider, without any loss of generality, that the world reference frame is aligned with the coordinate system of camera C_A , and that the X -axis is coincident with the baseline defined by the projection centers of the two pin-holes. The local reference frame of the second camera is assumed to have origin in C_B and to be parallel to the coordinate system of C_A . Under such circumstances the rigid transformation that maps point coordinates from C_B to C_A is given by

$$\mathbf{T}_{B \rightarrow A} = \begin{pmatrix} \mathbf{I} & \mathbf{h} \\ 0 & 1 \end{pmatrix} \quad (9)$$

with \mathbf{I} being the 3×3 identity matrix and $\mathbf{h} = (h \ 0 \ 0)^T$. Since the axes X, Y, Z of the system of coordinates of C_A , and the axes Y, Z of the reference frame of C_B are linearly independent lines, then they can be used to establish a basis G for the LLC defined by the baseline. It follows that

$$G \sim \begin{pmatrix} \mathbf{I} & \mathbf{I}^{\{2,3\}} \\ 0 & [\mathbf{h}]_{\times}^{\{2,3\}} \end{pmatrix} \quad (10)$$

with the upper script $\{2, 3\}$ denoting the second and third columns of the matrix.

Let us now consider an image correspondence $(\mathbf{x}, \mathbf{x}_A)$ between C and C_A . The back-projection of \mathbf{x}_A is a line \mathbf{L}_A with direction \mathbf{d}_A expressed in the reference frame of C_A . Given the basis G above, it follows that $\mathbf{L}_A \sim G \lambda_A$ with

$$\lambda_A \sim (\mathbf{d}_A^T \ 0 \ 0)^T. \quad (11)$$

Replacing in equation 8, and making $\mathbf{d} \sim \mathbf{K}^{-1}\mathbf{x}$, yields

$$\mathbf{x}^\top \mathbf{F}_A \mathbf{d}_A = 0 \quad (12)$$

with \mathbf{F}_A being the standard fundamental matrix between the uncalibrated camera \mathbf{C} and the calibrated view \mathbf{C}_A

$$\mathbf{F}_A = \mathbf{K}^{-\top} [\mathbf{t}]_\times \mathbf{R}. \quad (13)$$

Repeating the reasoning for the case of an image correspondence $(\mathbf{x}, \mathbf{x}_B)$ between \mathbf{C} and \mathbf{C}_B , it follows that

$$\mathbf{x}^\top \mathbf{F}_B \mathbf{d}_B = 0 \quad (14)$$

with \mathbf{F}_B being the fundamental matrix between \mathbf{C} and \mathbf{C}_B that can be written as

$$\mathbf{F}_B = \mathbf{F}_A + \mathbf{K}^{-1} \mathbf{R} [\mathbf{h}]_\times. \quad (15)$$

It follows from the equation above that the first columns of matrices \mathbf{F}_A and \mathbf{F}_B are always equal ($\mathbf{F}_A^{\{1\}} = \mathbf{F}_B^{\{1\}}$).

Given the image correspondences $(\mathbf{x}^{(i)}, \mathbf{d}_A^{(i)})$, with $i = 1, \dots, a$, and $(\mathbf{x}^{(a+j)}, \mathbf{d}_B^{(j)})$ with $j = 1, \dots, b$, we can determine the lines

$$\boldsymbol{\lambda}_A^{(i)} \sim \begin{pmatrix} d_{A,1}^{(i)} & d_{A,2}^{(i)} & d_{A,3}^{(i)} & 0 & 0 \end{pmatrix}^\top \quad (16)$$

$$\boldsymbol{\lambda}_B^{(i)} \sim \begin{pmatrix} d_{B,1}^{(i)} & 0 & 0 & d_{B,2}^{(i)} & d_{B,3}^{(i)} \end{pmatrix}^\top \quad (17)$$

and establish a system of linear equations based on the bilinear constraints of equations 12 and 14.

$$\begin{pmatrix} \boldsymbol{\lambda}_A^{(1)\top} \otimes \mathbf{x}^{(1)\top} \\ \vdots \\ \boldsymbol{\lambda}_A^{(a)\top} \otimes \mathbf{x}^{(a)\top} \\ \boldsymbol{\lambda}_B^{(1)\top} \otimes \mathbf{x}^{(a+1)\top} \\ \vdots \\ \boldsymbol{\lambda}_B^{(b)\top} \otimes \mathbf{x}^{(a+b)\top} \end{pmatrix} \begin{pmatrix} \mathbf{F}_A^{\{1\}} \\ \mathbf{F}_A^{\{2\}} \\ \mathbf{F}_A^{\{3\}} \\ \mathbf{F}_B^{\{2\}} \\ \mathbf{F}_B^{\{3\}} \end{pmatrix} = 0 \quad (18)$$

If $a + b \geq 14$ then the fundamental matrices \mathbf{F}_A and \mathbf{F}_B can be determined up to a common scale factor using a standard DLT approach.

5 MINIMAL SOLUTION

We have shown that the two fundamental matrices, \mathbf{F}_A and \mathbf{F}_B , that encode the calibration information \mathbf{K} , \mathbf{R} , and \mathbf{t} , can be determined from a minimum of 14 independent image correspondences. However, the total number of independent unknowns is 11 (5 intrinsic parameters and 6 extrinsic parameters) meaning that the estimation problem can be further constrained. Two of these constraints are rather obvious:

$$\det(\mathbf{F}_A) = 0, \quad (19)$$

$$\det(\mathbf{F}_B) = 0. \quad (20)$$

For the third constraint it must be observed that the sum of \mathbf{F}_A and \mathbf{F}_B is still a fundamental matrix.

From equations 13 and 15 it comes after algebraic manipulation that

$$\mathbf{F}_A + \mathbf{F}_B = \mathbf{K}^{-1} [2\mathbf{t} + \mathbf{R}\mathbf{h}]_\times \mathbf{R}, \quad (21)$$

which means that the following condition must hold

$$\det(\mathbf{F}_A + \mathbf{F}_B) = 0. \quad (22)$$

The equation above basically enforces the condition that \mathbf{F}_A and \mathbf{F}_B must be two fundamental matrices encoding the same rotation \mathbf{R} .

5.1 Outline of the estimation algorithm

\mathbf{F}_A and \mathbf{F}_B can be estimated from a minimum number of $a + b = 11$ pairwise correspondences. Note, however, that a single fundamental matrix can be estimated from 7 pairwise correspondences with a single camera, and therefore if $a > 7$ or $b > 7$ some equations are redundant. There are only two solvable minimal configurations in this problem: $(a = 7, b = 4)$ and $(a = 6, b = 5)$. We consider only the case $(a = 7, b = 4)$:

- 1) Build the linear system of equation 18 from the 11 pairwise correspondences.
- 2) Use the top 7 equations of this system determine a 2-dimensional solution space for the 9 parameters of \mathbf{F}_A using SVD. This enables to write $\mathbf{F}_A(\alpha) = \mathbf{A}' + \alpha \mathbf{A}$ with α being a free parameter.
- 3) Compute α by solving the cubic constraint of equation 19 and determine \mathbf{F}_A .
- 4) Substitute the up-to scale solution of \mathbf{F}_A in the linear system of equation 18. This system has now only 7 unknowns: the 6 parameters of \mathbf{F}_B and the scale factor of \mathbf{F}_A . The bottom 4 equations of this system can be used to determine a 3-dimensional solution space for \mathbf{F}_B . This enables to write $\mathbf{F}_B(\beta_1, \beta_2) = \mathbf{B}'' + \beta_1 \mathbf{B}' + \beta_2 \mathbf{B}$.
- 5) Substitute \mathbf{F}_A and $\mathbf{F}_B(\beta_1, \beta_2)$ in equations 20 and 22. This leads to a bivariate system of 2 quadratic equations. Compute β_1 and β_2 by solving the bivariate system [32], and determine the fundamental matrix \mathbf{F}_B .

Since the cubic equation of step 3 gives up to 3 discrete solutions, and the bivariate system of quadric equations has at most 4 distinct solutions, then there is a maximum of 12 possible solutions for the pair of fundamental matrices $(\mathbf{F}_A, \mathbf{F}_B)$.

5.2 Degenerate Configurations

The 11-point solution degenerates in two cases. If the 7 pairwise correspondences that are established with the same camera belong to a single plane the linear system of equation 18 is rank deficient. This is a known degenerate configuration in Fundamental matrix estimation []. The second degeneracy happens when there is no translation between the two calibrated cameras \mathbf{C}_A and \mathbf{C}_B . In this case all calibrated image rays belong to the same bundle and

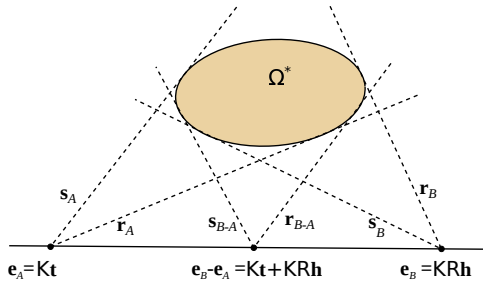


Fig. 4. Conic envelope Ω establishes linear relations $s^T K K^T s = 0$ and $r^T K K^T r = 0$.

the problem becomes equivalent to the estimation of a fundamental matrix between a calibrated and an uncalibrated pinhole views.

6 FACTORIZATION OF F_A AND F_B

In order to solve the calibration problem, F_A and F_B must be factorized into the intrinsic parameters K and the relative pose R, t . Let us first discuss the extraction of the matrix K . Consider the fundamental matrix F_A that is given in equation 13. After some algebraic manipulations we obtain that

$$F_A F_A^T \sim [e_A]_{\times} K K^T [e_A]_{\times} \quad (23)$$

with $e_A = Kt$ denoting the left side epipole of F_A (the image on C of the principal point of C_A). From the result above it follows that, if y is a point in the projective plane that satisfies

$$y^T F_A F_A^T y = 0, \quad (24)$$

then the line defined by y and e_A lies in the conic envelope $\Omega^* = K K^T$ that is the dual of the image of the absolute conic (DIAC) [12], [13]. $F_A F_A^T$ is a rank 2 symmetric matrix that can be understood as a degenerate conic locus comprising the points lying in two lines s_A, r_A that intersect e_A . From the two observations above it is easy to conclude that s_A, r_A must belong to the DIAC, as shown in Fig. 4. The same reasoning can be applied to the fundamental matrix F_B of equation 15

$$F_B F_B^T \sim [e_B]_{\times} K K^T [e_B]_{\times}, \quad (25)$$

and to the matrix $F_B - F_A$ that is still rank deficient because the first columns of the two fundamental matrices are equal

$$(F_B - F_A)(F_B - F_A)^T \sim [e_B - e_A]_{\times} K K^T [e_B - e_A]_{\times} \quad (26)$$

Summarizing, and as shown in Fig. 4, the DIAC is fully constrained by the line pairs arising from the rank 2 degenerate conics $F_A F_A^T, F_B F_B^T$, and $(F_B - F_A)(F_B - F_A)^T$. It is important to refer that, although we have six lines, they only give rise to five independent constraints on the parameters of the DIAC. This is explained by the fact that their pairwise intersections are collinear.

After knowing K , we can compute the essential matrix E_A and apply standard techniques for determining the rotation R and the translation t up to scale factor [12], [13]. The scale factor can be easily found using the known baseline between C_A and C_B .

7 RANSAC WITH MULTIPLE DATASETS

RANSAC [15] is the most widely used method to eliminate outlier correspondences when a minimal solution is available. This method attempts to fit a model T to a single dataset \mathcal{D} with L correspondences which are either inliers or outliers. RANSAC iteratively generates candidate models $T_C = g(\mathcal{S})$ by randomly selecting a subset \mathcal{S} with s random samples from \mathcal{D} (Fig. 5(a)). In each iteration a candidate model T_C is evaluated using some cost metric and whenever a model with a lower cost is found it is stored as the current best candidate. After a certain number of iterations n RANSAC stops and outputs the best candidate. Different versions of RANSAC have different cost metrics: original RANSAC [15] minimizes the number of outliers for a pre-defined threshold t ; MLESAC [22] chooses the model with maximum likelihood, assuming that the residue of inliers follows a gaussian distribution, while the residue of outliers follow a uniform distribution; MAPSAC [23] maximizes the posterior probability of a model and its latent parameters. Despite these differences, all versions of RANSAC work under the assumption that samples are selected from a single dataset \mathcal{D} with a certain inlier ratio γ whose value is updated according to the current best candidate. An accurate estimation of the inlier ratio γ is important to know the required number of RANSAC iterations and also to compute the cost metrics of MLESAC and MAPSAC.

Our problem, however, does not fit into the standard assumptions of RANSAC. A model generator for our problem involves selecting 7 correspondences from one dataset and 4 from another. These two datasets might have different inlier ratios and thus all RANSAC assumptions that depend on a single value γ must be revised. Additionally we might think of a scenario where there are correspondences with $N > 2$ cameras and thus to use RANSAC with our algorithm we must first select 2 cameras and only then sample 7 and 4 correspondences from them. With these issues in mind we propose a new framework, called multi-RANSAC, that takes into account the sampling of different datasets.

For the sake of generalization we assume an arbitrary problem with N datasets $\mathcal{D}_1, \dots, \mathcal{D}_N$ and a model generator $T_C = g(\mathcal{S}_1, \dots, \mathcal{S}_M)$ that requires M subsets \mathcal{S}_j , each of them containing s_j samples from one of the datasets $\mathcal{D}_1, \dots, \mathcal{D}_N$. As displayed in Fig. 5(b), the sampling process is done in two steps: it first randomly selects M datasets $\hat{\mathcal{D}}_1, \dots, \hat{\mathcal{D}}_M$ from the N datasets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_N$, with $M \leq N$; then it randomly

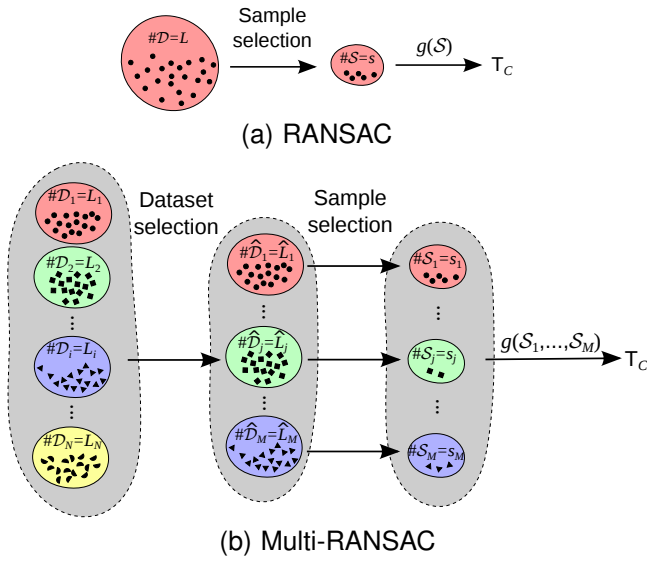


Fig. 5. Generation of candidate models T_C . (a) In each RANSAC iteration a subset S with s samples is selected from dataset D . (b) In each multi-RANSAC iteration there are two sampling steps: the first step randomly selects M of the N datasets D_1, \dots, D_N ; in the second step s_j samples are selected from each of the M selected datasets \hat{D}_j .

selects a subset S_j with s_j samples from each selected dataset \hat{D}_j .

In this section we discuss the necessary adaptations to RANSAC, MLESAC, and MAPSAC when dealing with multiple datasets, which we designate by multi-RANSAC, multi-MLESAC, and multi-MAPSAC respectively.

7.1 Multi-RANSAC

In standard RANSAC it is assumed that inlier samples have a uniform error distribution over some bounded interval. All samples with an error greater than a threshold t are considered outliers. In this case the evaluation cost of each candidate model is simply the total number of outliers. For the multi-RANSAC approach we can use the same evaluation metric by summing up the outliers in all datasets. We might also consider tuning different thresholds for each dataset, when it makes sense in a particular problem.

In standard RANSAC the number of iterations n is determined by guaranteeing that at least in one of the iterations a model is generated only from inlier samples with a probability p , set to a value close to 1. The sampling process is approximated by a succession of s Bernoulli trials, i. e., a succession of s independent sample selections with a constant probability γ of selecting an inlier in each of them. Therefore, the probability p_{ins} of selecting s consecutive inliers is

$$p_{ins} = \gamma^s \quad (27)$$

Note that the probability γ also represents the inlier ratio in dataset D . We want to guarantee that the probability of never selecting s inliers after n iterations is lower than $1 - p$, i. e.

$$(1 - \gamma^s)^n < 1 - p. \quad (28)$$

Therefore, the number of RANSAC iterations is

$$n = \frac{\log(1 - p)}{\log(1 - \gamma^s)}. \quad (29)$$

Whenever a new best model is found, the values γ and consequently n are updated.

In the multi-RANSAC case n must be computed differently since the sample selection process has two steps and each dataset D_i might have a different inlier ratio γ_i . The probability of obtaining an inlier by first selecting a random dataset D_i and then selecting a random sample from it is

$$p_{in} = \frac{1}{N} \sum_{i=1}^N \gamma_i. \quad (30)$$

In an analogous manner to the standard RANSAC formulation we assume that the probability of selecting an inlier from dataset D_i is a constant value γ_i for successive selections (i. e. the second selection step is a succession of Bernoulli trials). The probability of selecting s_j inliers by first selecting a random dataset D_i and then selecting s_j random samples from it is

$$p_{ins} = \frac{1}{N} \sum_{i=1}^N \gamma_i^{s_j}. \quad (31)$$

We now further approximate the first selection step (dataset selection) by a succession of Bernoulli trials. The complete multi-RANSAC sampling process is thus simplified to selecting a random dataset M times and for each of them successively selecting s_1, s_2, \dots, s_M samples. The probability of selecting only inliers in this process is

$$p_{inall} = \prod_{j=1}^M \frac{1}{N} \sum_{i=1}^N \gamma_i^{s_j}. \quad (32)$$

Using the same reasoning behind equation 29, the number of multi-RANSAC iterations is now

$$n = \frac{\log(1 - p)}{\log(1 - p_{inall})}. \quad (33)$$

Note however that approximating the dataset selection step by Bernoulli trials is only valid when M is much smaller than N because the same dataset cannot be selected more than once. Specifically in the case that $M = N$ the above equation might grossly misestimate the number of required iterations. In this case the dataset selection step outputs a random permutation of all datasets, choosing for each dataset D_i which number of samples s_j is selected. There are $N!$ possible dataset selections, and to obtain p_{inall} we

must weigh in the probability of selecting only inliers for all possible permutations.

To illustrate this case consider our calibration problem when there are only correspondences with two cameras. In this case we have two datasets with pairwise correspondences $\mathcal{D}_A, \mathcal{D}_B$ with inlier ratios γ_A, γ_B and we want to select 7 correspondences from one of them and 4 from the other. In this case $N = M = 2$ and equation 32 becomes

$$p_{inall} = \frac{1}{4} (\gamma_A^7 \gamma_B^4 + \gamma_B^7 \gamma_A^4 + \gamma_A^7 \gamma_A^4 + \gamma_B^7 \gamma_B^4) \quad (34)$$

However, in a more careful analysis, we can observe that the dataset selection step has only two possible outcomes: either 7 correspondences are selected from \mathcal{D}_A and 4 from \mathcal{D}_B , or 4 correspondences are selected from \mathcal{D}_A and 7 from \mathcal{D}_B . The probability of selecting only inliers in this process is

$$p_{inall} = \frac{1}{2} (\gamma_A^7 \gamma_B^4 + \gamma_B^7 \gamma_A^4). \quad (35)$$

All the results derived for computing the number of multi-RANSAC iterations also extend to the multi-MLESAC and multi-MAPSAC formulations discussed next.

7.2 Multi-MLESAC

MLESAC [22] aims at finding the model T with minimum negative log-likelihood, given a set of measurements \mathcal{D} . Each sample \mathbf{d}_k in \mathcal{D} can be put into one of two subsets: the inliers \mathcal{I} or the outliers \mathcal{O} .

The residue of samples in \mathcal{I} is assumed to follow a gaussian distribution $N(0, \sigma)$. A model T , given an inlier sample \mathbf{d}_k with residue $r_k^{\mathcal{I}}$, has a likelihood

$$L(T|r_k^{\mathcal{I}}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{|r_k^{\mathcal{I}}|^2}{2\sigma^2}}. \quad (36)$$

The samples from \mathcal{O} are observations independent from the model, and their residue is assumed to follow a uniform distribution over an interval $[-\frac{v}{2}, \frac{v}{2}]$. A model T , given an outlier sample \mathbf{d}_k with residue $r_k^{\mathcal{O}}$ has a constant likelihood

$$L(T|r_k^{\mathcal{O}}) = \frac{1}{v}. \quad (37)$$

The samples from dataset \mathcal{D} follow a mixed distribution of inliers and outliers and therefore the likelihood $L(T|r_k)$ of a model T , given a random sample \mathbf{d}_k from \mathcal{D} with residue r_k is

$$L(T|r_k) = \left(\gamma \left(\frac{1}{\sqrt{2\pi}\sigma} \right) e^{-\frac{|r_k|^2}{2\sigma^2}} + \frac{1-\gamma}{v} \right), \quad (38)$$

where γ is the probability of \mathbf{d}_k being an inlier.

The MLESAC problem can now be formulated by considering the negative log-likelihood of T given all L samples in \mathcal{D}

$$\min_T \left(- \sum_{k=1}^L \log L(T|r_k) \right). \quad (39)$$

Note that the inlier ratio γ is updated in each iteration using expectation maximization with the following constraint:

$$\gamma = \frac{1}{L} \sum_{k=1}^L Pr(r_k^{\mathcal{I}}|\gamma) \quad (40)$$

where γ is initialized to 0.5 on the left side of the equation and is iteratively updated until convergence.

We now consider the multi-MLESAC problem. When sampling from N different datasets we aim at maximizing the likelihood of model T given datasets $\mathcal{D}_1, \dots, \mathcal{D}_N$, each of them with a number of samples L_i , an inlier standard deviation σ_i , an outlier range v_i , and an inlier ratio γ_i . In this case the likelihood of a model T , given a sample $\mathbf{d}_{i,k}$ from dataset \mathcal{D}_i with a residue $r_{i,k}$ is

$$L(T|r_{i,k}) = \left(\gamma_i \left(\frac{1}{\sqrt{2\pi}\sigma_i} \right) e^{-\frac{|r_{i,k}|^2}{2\sigma_i^2}} + \frac{1-\gamma_i}{v_i} \right). \quad (41)$$

The multi-MLESAC problem for N datasets can now be formulated as

$$\min_T \left(- \sum_{i=1}^N \sum_{k=1}^{L_i} \log L(T|r_{i,k}) \right). \quad (42)$$

Note that to compute $\gamma_1, \dots, \gamma_N$ in each iteration we have to solve N expectation maximization problems with the form of equation 40.

After multi-MLESAC is finished, the inliers of the best candidate model can be found by checking for each sample if its probability of being an inlier is higher than of being an outlier

$$\gamma_i L(T|r_{i,k}^{\mathcal{I}}) > (1-\gamma_i) L(T|r_{i,k}^{\mathcal{O}}) \quad (43)$$

which, by observation of equations 36 and 37, can be rewritten as

$$|r_{i,k}|^2 < -2\sigma_i^2 \ln \frac{\sqrt{2\pi}\sigma_i(1-\gamma_i)}{\gamma_i v_i}. \quad (44)$$

The most notable difference when we step from a standard MLESAC formulation to multi-MLESAC is that different datasets might have different inlier ratios γ_i . This reflects a practical scenario where some datasets are consistently more reliable than others. Multi-MLESAC is able to capture those differences by estimating separate values γ_i for each dataset, which in turn results in a different cost function and inlier threshold for each dataset.

7.3 Multi-MAPSAC

The MLESAC formulation can be further generalized to a maximum a posteriori problem (MAPSAC [23]). While [23] does a very exhaustive bayesian analysis of random sampling for geometric problems, we are only interested in its key observation that an algorithm from the RANSAC family does not only estimate the parameters of model T but also an additional set

of latent parameters, namely by deciding whether each sample is an inlier or an outlier through the expectation maximization of the inlier ratio γ . Taking this into account we formulate the multi-MAPSAC problem as

$$\max_{\mathbf{T}, \gamma_1, \dots, \gamma_N} Pr(\mathbf{T}, \gamma_1, \dots, \gamma_N | \mathbf{R}_1, \dots, \mathbf{R}_N), \quad (45)$$

where $\mathbf{R}_1, \dots, \mathbf{R}_N$ represent the residues of all samples from $\mathcal{D}_1, \dots, \mathcal{D}_N$ respectively, which follow the mixed inlier-outlier distribution of multi-MLESAC. This formulation can be re-written as

$$\max_{\mathbf{T}, \gamma_1, \dots, \gamma_N} Pr(\gamma_1, \dots, \gamma_N, \mathbf{T}) Pr(\mathbf{R}_1, \dots, \mathbf{R}_N | \mathbf{T}, \gamma_1, \dots, \gamma_N) \quad (46)$$

Note that although this is a MAP formulation, it is not a step-by-step generalization of the MAPSAC method described in [23], which deals with the marginalization of parameter γ and the effect of additional latent parameters, e. g. reconstructed 3D points. In this paper we do not take these issues into account.

When compared to multi-MLESAC, equation 46 has an additional prior on the model and the latent parameters $Pr(\gamma_1, \dots, \gamma_N, \mathbf{T})$. Prior knowledge about the model \mathbf{T} is a very specific issue in each application scenario and we ignore it in the context of this paper. Our main motivation behind this formulation is to account for prior knowledge about the inlier ratios γ_i . While multi-MLESAC assumes that parameters γ_i are independent from each other, with multi-MAPSAC we want to account for the possibility that this is not the case.

Using prior knowledge on the relative distribution of inlier ratios γ_i is important in the context of our calibration problem. For simplicity, consider the case where there are pairwise correspondences with only two cameras ($N = M = 2$). Correspondences with one camera just give us a fundamental matrix, while the pairwise correspondences with two cameras give us both the extrinsic and intrinsic camera calibration. This means that a candidate model with many inliers in one dataset and very few on the other is a poor solution that is over-fitting to a particular fundamental matrix. To tackle this issue we use the multi-MAPSAC approach and define a prior probability function $Pr(\gamma_A, \gamma_B)$ that penalizes significantly uneven distributions of inliers

$$Pr(\gamma_A, \gamma_B) = (\alpha + 1)^2 (\gamma_A \gamma_B)^\alpha, \quad (47)$$

where the parameter α is set to a value with the same order of magnitude as the number of correspondences in each dataset. Note that the constant factor $(\alpha + 1)^2$ is just to guarantee that $Pr(\gamma_A, \gamma_B)$ is a probability density function for γ_A, γ_B between 0 and 1. In the context of maximum a posteriori estimation it can be ignored. These observations also extend to a scenario where there are correspondences with N cameras C_1, \dots, C_N . Note, however, that the over-fitting case

discussed above can only happen when a candidate solution fits well into just one dataset. Thus, the prior term $Pr(\gamma_A, \gamma_B)$ should be computed using only the two highest values from $\gamma_1, \dots, \gamma_N$.

8 CAMERA CALIBRATION WITH N VIEWS

The multi-MAPSAC formulation can be used together with the minimal solution described in section 5 to provide a generalised calibration pipeline that is able to add a new node into a network given an arbitrary number N of different calibrated views. To achieve this, in each multi-MAPSAC iteration, we start by randomly sampling two out of N calibrated views and then sample 7 and 4 correspondences from each of those views respectively. Each candidate solution is then evaluated against all available datasets according to equation 46. It must be noted, however, that in a realistic scenario a large camera network will contain many different nodes that do not overlap with the field of view of the newly added camera. This would make the multi-MAPSAC sampling process extremely inefficient, since in many iterations datasets with no inliers would be selected. To avoid this we perform a pre-filtering step where only the views with enough correspondences are considered for the sampling process. The complete calibration algorithm is summarised in the following section.

8.1 Algorithm Outline

Given a set of N calibrated views, a new node can be fully calibrated with the following process:

- 1) Establish pairwise point correspondences between the new node and all the N calibrated views (e. g. using SIFT features)
- 2) For each of the N sets of pairwise correspondences, perform 7-point RANSAC and remove outlier correspondences
- 3) Out of the N sets of pairwise correspondences, select those whose number of correspondences surpasses a given threshold
- 4) Perform multi-RANSAC on the selected sets, providing an initial calibration solution and a set of inlier pairwise correspondences
- 5) Refine the solution using Bundle Adjustment, as detailed in section 8.2.

8.2 Bundle adjustment

A final refinement with bundle adjustment must be used to achieve an optimal solution. Usually bundle adjustment minimizes the re-projection of reconstructed 3D points onto the cameras. However, since our formulation only uses pairwise correspondences, the introduction of unknown 3D points is an unnecessary burden. As described in [13], an explicit representation of 3D points can be avoided by minimizing the perpendicular distances between point correspondences and their epipolar lines.

Given a pairwise point correspondence $(\mathbf{x}, \hat{\mathbf{x}})$ between two cameras related by a fundamental matrix F , the epipolar error r can be measured by the distance in pixels between point \mathbf{x} and the epipolar line of $\hat{\mathbf{x}}$

$$r = \frac{\mathbf{x}^\top F^\top \hat{\mathbf{x}}}{\|\mathbb{I}_{2 \times 3} F^\top \hat{\mathbf{x}}\|}. \quad (48)$$

with

$$\mathbb{I}_{2 \times 3} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}. \quad (49)$$

Analogously, the distance between point $\hat{\mathbf{x}}$ and the epipolar line of \mathbf{x} is

$$\hat{r} = \frac{\hat{\mathbf{x}}^\top F \mathbf{x}}{\|\mathbb{I}_{2 \times 3} F \mathbf{x}\|}. \quad (50)$$

We now consider a network with N calibrated cameras with rotations $\{R_1, R_2, \dots, R_N\}$, translations $\{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_N\}$, and intrinsics $\{K_1, K_2, \dots, K_N\}$ in a common reference frame, and a new camera with unknown parameters R, \mathbf{t}, K . The new camera has a set of L_i pairwise correspondences $\{(\mathbf{x}_{i,1}, \hat{\mathbf{x}}_{i,1}), (\mathbf{x}_{i,2}, \hat{\mathbf{x}}_{i,2}), \dots, (\mathbf{x}_{i,L_i}, \hat{\mathbf{x}}_{i,L_i})\}$ with each calibrated camera $i = 1, 2, \dots, N$. Therefore, the bundle adjustment problem becomes

$$\min_{R, \mathbf{t}, K} \sum_{i=1}^N \sum_{k=1}^{L_i} r_{i,k}^2 + \hat{r}_{i,k}^2 \quad (51)$$

with

$$r_{i,k} = \left(\frac{\mathbf{x}_{i,k}^\top F_i^\top \hat{\mathbf{x}}_{i,k}}{\|\mathbb{I}_{2 \times 3} F_i^\top \hat{\mathbf{x}}_{i,k}\|} \right) \quad (52)$$

$$\hat{r}_{i,k} = \left(\frac{\hat{\mathbf{x}}_{i,k}^\top F_i \mathbf{x}_{i,k}}{\|\mathbb{I}_{2 \times 3} F_i \mathbf{x}_{i,k}\|} \right) \quad (53)$$

$$F_i = K_i [R_i^\top \mathbf{t} + \mathbf{t}_i]_\times R_i^\top R K. \quad (54)$$

Although we do not consider lens distortion in this paper, it could be incorporated into the bundle adjustment cost function by considering the epipolar error $r_{i,k}$ for a correspondence $(\mathbf{x}, \hat{\mathbf{x}})$ as the minimum distance between a point \mathbf{x} and the epipolar curve defined by $\hat{\mathbf{x}}, F$, and the distortion model.

9 EXPERIMENTS

In this section we validate our calibration method using both synthetic data and real imagery. In a first set of experiments we use synthetic data to demonstrate challenging scenarios where our multi-RANSAC formulation is essential to obtain accurate calibrations. We then demonstrate the usefulness of our calibration method in practice with three distinct experiments. The first experiment concerns adding, or re-calibrating, a camera node during network operation, using image point correspondences at a certain frame time instant. The second experiment refers to recovering the trajectory and intrinsics of a free moving camera whose acquisition is synchronized with

the network. In the last experiment we demonstrate the usefulness of our algorithm in the context of Structure-from-Motion [33], as a post-processing tool to reconstruct challenging viewpoints. In all experiments, we use SIFT features to establish point correspondences between views.

In the experiments with real data we compare our 11-point calibration algorithm as described in section 8.1 against the standard approach for calibrating a new view from 3D points in the scene, i. e. the 6-point linear algorithm [12]. The former uses independent pairwise correspondences with two calibrated views, while the latter requires triple correspondences such that each point in the uncalibrated images is seen by at minimum of two calibrated cameras in order to enable 3D reconstruction. While our method requires using multi-MAPSAC and bundle adjustment as described in section 8.1, the 6-point approach is a single dataset formulation and relies on 3D point estimation, therefore we use the standard versions of both MLESAC and bundle adjustment.

9.1 Validation of multi-MAPSAC

We built a simulated environment in order to investigate the conditions where a multi-dataset RANSAC approach is fundamental to obtain reliable calibration results. In each simulation trial we generate calibrated cameras C_A, C_B , and an uncalibrated camera C in random poses within a spherical volume such that they share a common field of view. Then we randomly generate 3D points that are viewed by cameras C_A, C_B , and C . All correspondences are injected with gaussian noise with 1 pixel standard deviation and a predefined ratio of outliers. We tried to calibrate camera C using the 11-point algorithm with multi-MAPSAC, multi-MLESAC and standard MLESAC for different distributions in the number of correspondences and outlier ratios. For a varying number of correspondences we consider 10% outliers in C_A and 40% outliers in C_B , while for a varying outlier ratio we consider 500 correspondences with C_A and 100 correspondences with C_B . It can be observed that our multi-MLESAC and multi-MAPSAC approaches outperform standard MLESAC when there is an imbalance in both outlier ratio and number of correspondences. Additionally, when this imbalance is too great the use of a prior on the outlier distribution becomes important to avoid overfitting to C_A , and thus multi-MAPSAC significantly outperforms multi-MLESAC. Note that a simultaneous imbalance in number of correspondences and outlier ratio is a scenario that naturally occurs in practice, since with larger baselines and greater changes between two viewpoints, not only feature matching methods generally detect fewer correspondences, but the number of mismatches and the ratio of outliers increases as well.

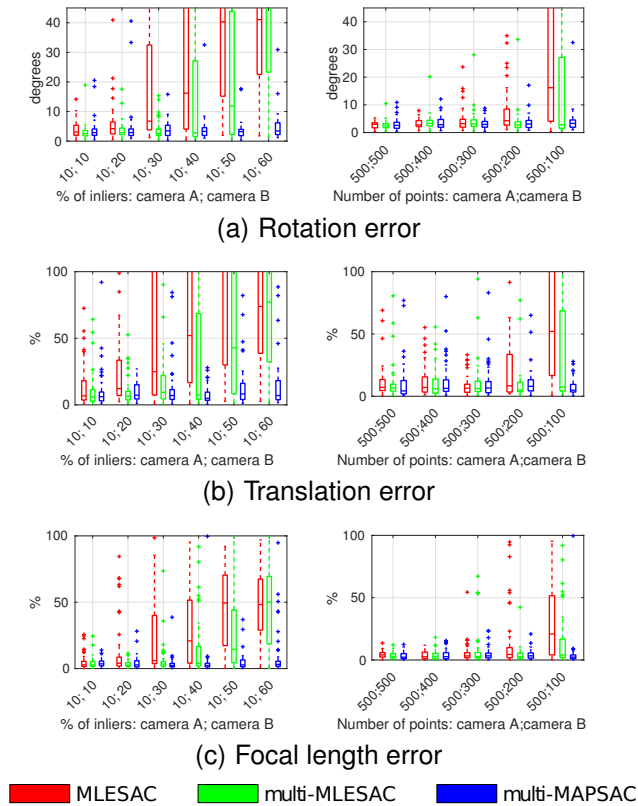
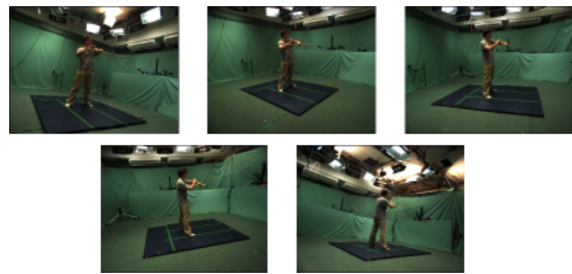


Fig. 6. Comparison between multi-MAPSAC, multi-MLESAC and MLESAC using synthetic data. Each error distribution contains results from 50 calibration trials. The results with varying number of correspondences (a, c, e) contain 10% outliers in camera A and 40% outliers in camera B. The results with varying number of correspondences (b, d, f) contain 500 correspondences in camera A and 100 correspondences in camera B.

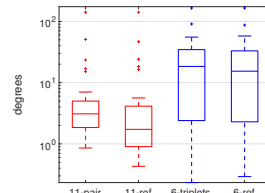
9.2 Addition of a new node to a calibrated network

In this experiment we aim at fully calibrating a camera using pairwise correspondences with a set of frames acquired at the same time instant. Data was acquired with the Grimage platform [34] that comprises a set of synchronized camera nodes in a room, and it is publicly available as the *stick* dataset in the 4drepository [35]. The cameras are calibrated both intrinsically and extrinsically with the method described in [9].

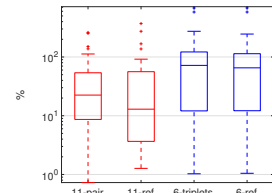
We want to compare our 11-point approach against the 6-point method using all available correspondences across the 5 views. For this purpose we selected a particular time instant of the dataset and the five views that are shown in Fig. 7(a). We then tried to calibrate each of them assuming the remaining four were calibrated. The selected camera nodes in this experiment have significant changes in viewpoint, making it very difficult to establish triple correspondences. We want to show that in many situations there are enough pairwise correspondences for our 11-point algorithm to provide accurate results but not enough



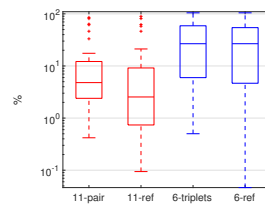
(a) Input images



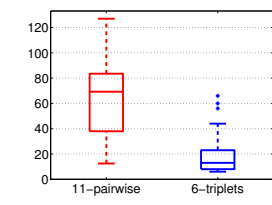
(b) Rotation error



(c) Translation error



(d) Focal length error



(e) Number of inliers

Fig. 7. Addition of a new node to a camera network. In each trial we try to calibrate one of the cameras in (a) assuming that the remaining four are calibrated. (b), (c), (d), (e) show the comparative performance between 11-point (pairwise) and 6-point (triple) for 250 calibration trials.

triple correspondences for the 6-point algorithm to work.

Since there is a wide baseline between the five cameras the two closest cameras typically produce the majority (if not all) the reliable correspondences. Therefore, the pipeline described in 8.1 always select the two views with the higher number of correspondences. In the case of the 6-point algorithm all pre-filtered triple correspondences from the four cameras are used.

For each of the five cameras, the calibration with the two methods is carried 50 times, summing up to 250 calibration tests for each approach. The error distributions before (11-pair, 6-triple) and after (11-ref, 6-ref) the bundle adjustment step are provided in Fig. 7. Note that the errors are displayed in logarithmic scale and our algorithm provides extremely more accurate results than the 6-triplets approach, which completely fails to provide a reasonable calibration in most cases. This can be explained by the fact that it is possible to establish a much higher number of pairwise correspondences than triple correspondences (Fig. 7(e)), despite the fact that triple correspondences

are established across the four calibrated cameras.

9.3 Calibration of a hand-held camera

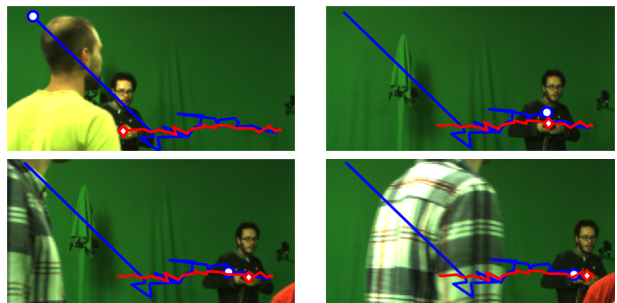
Using the same platform of the previous experiment, we acquired a set of video sequences where one of the synchronized cameras is a hand-held moving device. Each sequence is composed of 30 frames in which the hand-held camera shares its field of view with two other calibrated cameras (Fig. 8). In comparison with the previous experiment the calibrated camera nodes have a smaller baseline and the viewed scene is richer in features. This benefits the standard 6-point approach, as it is easier to establish triple correspondences. However, the viewed scene is highly dynamic and contains significant occlusions in some frames, making it difficult to establish triple correspondences. The intrinsic parameters of the hand-held camera were previously determined using the method described in [9] and we use these values as groundtruth for comparison with our estimates.

Both the intrinsic parameters and the trajectory of the hand-held camera are recovered with both our 11-point method and the 6-point method. In a first step, we calibrate each frame independently using pairwise correspondences with the synchronized frames from the calibrated cameras. This is convenient for the case of a hand-held camera with motorized lenses for which the zoom varies while moving. However, since in this experiment we know that the camera intrinsics are stationary, a final estimation with bundle adjustment is made assuming a single set of intrinsics for all frames.

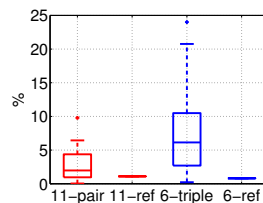
The error distribution for the intrinsic parameters before and after global refinement is presented in Fig. 8(c) and 8(d). Note that although the results for the standard 6-point approach are worse than our 11-point approach, they are significantly better than in the previous experiment. As explained earlier this is to be expected, since it is easier to establish correspondences in this set of acquisitions. Although the initialization results are sufficient for the focal length to converge to similarly accurate values with both algorithms, our algorithm is able to provide a much better estimation of the principal point. Our estimated camera trajectory is also significantly smoother and in line with a reasonable hand-held trajectory (Fig. 8(f)), specially when significant occlusions occur (e. g. the leftmost frames in Fig. 8(e)). Since we do not have groundtruth values for the camera trajectory, in Fig. 8(b) both trajectories are projected onto the image plane of a third calibrated camera in which the person handling the free camera is visible. For the selected frames it is quite clear that our estimations with the 11-point approach (red) are significantly more accurate than with the 6-point approach (blue). This confirms the intuition from Fig. 8(f) that our algorithm provides accurate trajectory estimations.



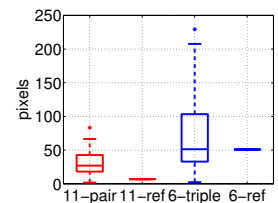
(a) Sample Frames



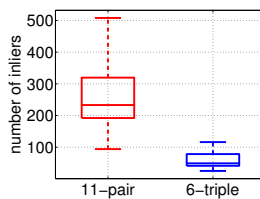
(b) Projected Trajectory



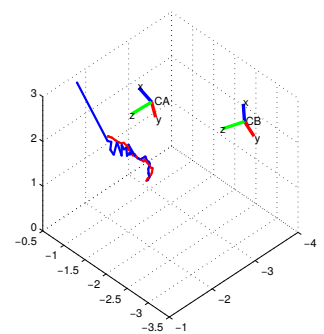
(c) Focal length



(d) Principal point



(e) Number of inliers



(f) Camera trajectory

Fig. 8. Online calibration of a hand-held camera. Sample images from the synchronized video sequence are depicted in (a). The trajectory of the hand-held camera was estimated with the 11-point (blue) and the 6-point (red) algorithms. In (b) the estimated trajectories are projected on a different calibrated view.

9.4 Addition of new nodes to an SfM reconstruction

While we mainly focus on camera networks, our calibration algorithm can be used in any application that involves calibrating a camera using multiple cali-

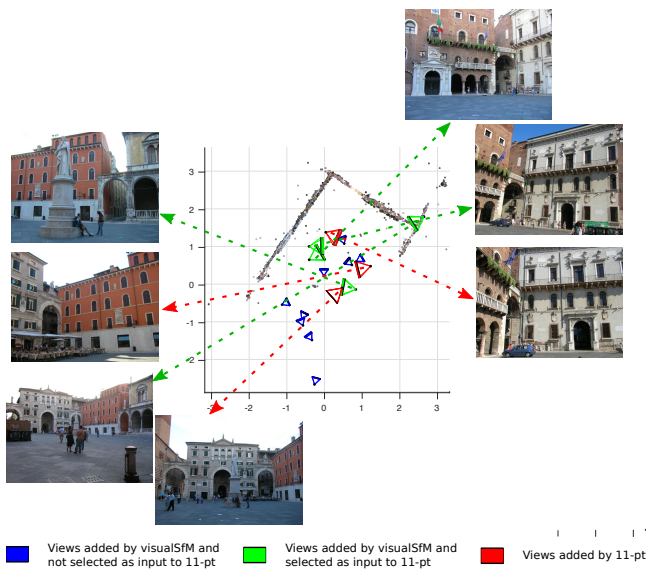


Fig. 9. VisualSfM was run on a set of 19 uncalibrated images and generated a 3D model with 13 views (blue and green). From the remaining 6, our 11-point calibration algorithm was able to insert 3 more views into the SfM model (red).

brated viewpoints, such as a SfM pipeline [33]. Traditional SfM works by calibrating an initial pair of views and incrementally adding a new image to the model in each iteration. New viewpoints are initialised using correspondences with 3D points in the model (i. e. triple correspondences) and then refined with bundle adjustment. SfM stops when it is unable to find any more images with enough triple correspondences.

In this section we show that some of the viewpoints that an SfM pipeline is unable to reconstruct due to the lack of triple correspondences can still be recovered using pairwise correspondences. We ran visualSfM [36] on a collection of 19 images from the publicly available *Piazza Dante* dataset [37]. The output reconstruction is displayed in Fig. 9, where the blue and green markers represent the 13 viewpoints that visualSfM is able to reconstruct. On the remaining 6 viewpoints that visualSfM is unable to reconstruct, our calibration pipeline is able to detect 3 views (red markers in Fig. 9) that have more than 50 pairwise correspondences with at least two already reconstructed views (green markers in Fig. 9). Note that the viewpoints captured by our method are mainly observing regions that do not contain 3D reconstructed points. In this way our 11-point algorithm can potentially be used as an additional tool within SfM to insert very challenging views.

10 CONCLUSION

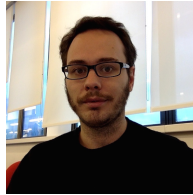
We presented a new method for the intrinsic and extrinsic calibration of a camera from pairwise correspondences with other calibrated cameras. We

showed that there are practical scenarios in the context of camera networks and Structure-from-Motion where large baselines and limited overlap between viewpoints makes it extremely difficult or impossible to use traditional methods that rely on 3D points to perform camera calibration. In these challenging cases our method offers a new alternative solution that solely relies on pairwise correspondences. One of the key aspects of our method is the use of multi-RANSAC, which is able to deal with the sampling of pairwise correspondences across multiple views taking into account that they might induce different error distributions. This framework is essential for our 11-point minimal solution to work in practice. There are many aspects of the multi-RANSAC framework that can be further developed. While in this paper we mostly focus on analysing the effect of different inlier distributions, other factors such as differences in the error distribution of inlier samples might play an important role in other multi-RANSAC applications. This framework is also potentially useful in other problems involving multiple views or different types of correspondences.

REFERENCES

- [1] J. Starck and A. Hilton, "Surface capture for performance-based animation," *Computer Graphics and Applications, IEEE*, vol. 27, no. 3, pp. 21–31, 2007.
- [2] Z. Zhao and Y. Liu, "Practical multi-camera calibration algorithm with 1d objects for virtual environments," in *Multimedia and Expo, 2008 IEEE International Conference on*, 2008, pp. 1197–1200.
- [3] E. Shen and R. Hornsey, "Multi-camera network calibration with a non-planar target," *Sensors Journal, IEEE*, vol. 11, no. 10, pp. 2356–2364, 2011.
- [4] J. Courchay, A. Dalalyan, R. Keriven, and P. Sturm, "A global camera network calibration method with linear programming," in *Proceedings of the International Symposium on 3D Data Processing, Visualization and Transmission*, 2010.
- [5] R. Kumar, A. Ilie, J.-M. Frahm, and M. Pollefeys, "Simple calibration of non-overlapping cameras with a mirror," in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, June 2008, pp. 1–7.
- [6] R. Rodrigues, J. a. P. Barreto, and U. Nunes, "Camera pose estimation using images of planar mirror reflections," in *Computer Vision ECCV 2010*, ser. Lecture Notes in Computer Science, K. Daniilidis, P. Maragos, and N. Paragios, Eds. Springer Berlin Heidelberg, 2010, vol. 6314, pp. 382–395.
- [7] T. Svoboda, D. Martinec, and T. Pajdla, "A convenient multicamera self-calibration for virtual environments," *Presence: Teleoper. Virtual Environ.*, vol. 14, no. 4, pp. 407–422, 2005.
- [8] J. Barreto and K. Daniilidis, "Wide area multiple camera calibration and estimation of radial distortion," in *OMNIVIS'2004 - Int. Workshop in Omnidirectional vision, camera networks, and non-conventional cameras*, 2004.
- [9] A. Zaharescu, R. Horaud, R. Ronfard, and L. Lefort, "Multiple camera calibration using robust perspective factorization," in *3D Data Processing, Visualization, and Transmission, Third International Symposium on*, 2006, pp. 504–511.
- [10] G. Thomas, "Real-time camera tracking using sports pitch markings," *Journal of Real-Time Image Processing*, vol. 2, no. 2-3, pp. 117–132, 2007.
- [11] J. Puwein, R. Ziegler, J. Vogel, and M. Pollefeys, "Robust multi-view camera calibration for wide-baseline camera networks," in *Applications of Computer Vision (WACV), 2011 IEEE Workshop on*, Jan 2011, pp. 321–328.
- [12] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge Academic Press, 2003.

- [13] Y. Ma, S. Soatto, J. Kosecka, and S. Sastry, *An invitation to 3-D vision: from images to geometric models*. Springer, 2004.
- [14] F. Vasconcelos, J. P. Barreto, and E. Boyer, "A minimal solution for camera calibration using independent pairwise correspondences," in *Computer Vision—ECCV 2012*. Springer, 2012, pp. 724–737.
- [15] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [16] B. Clipp, C. Zach, J.-M. Frahm, and M. Pollefeys, "A new minimal solution to the relative pose of a calibrated stereo camera with small field of view overlap," in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 1725–1732.
- [17] F. Vasconcelos and J. P. Barreto, "Towards a minimal solution for the relative pose between axial cameras," in *BMVC*, 2013.
- [18] R. Pless, "Using many cameras as one," in *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, 2003.
- [19] H. Stewénius, D. Nistér, M. Oskarsson, and K. Åström, "Solutions to minimal generalized relative pose problems," in *Workshop on Omnidirectional Vision*, Beijing China, 2005.
- [20] C. Raposo, M. Lourenço, J. P. Barreto, and M. Antunes, "Plane-based odometry using an rgb-d camera," in *BMVC*, 2013.
- [21] C. Raposo, M. Antunes, and J. P. Barreto, "Piecewise-planar stereoscan: structure and motion from plane primitives," in *Computer Vision – ECCV 2014*, ser. Lecture Notes in Computer Science, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Springer International Publishing, 2014, vol. 8690, pp. 48–63.
- [22] P. H. Torr and A. Zisserman, "Mlesac: A new robust estimator with application to estimating image geometry," *Computer Vision and Image Understanding*, vol. 78, no. 1, pp. 138–156, 2000.
- [23] P. H. S. Torr, "Bayesian model estimation and selection for epipolar geometry and generic manifold fitting," *International Journal of Computer Vision*, vol. 50, no. 1, pp. 35–61, 2002.
- [24] J. Kim, L. Hodong, and R. Hartley, "Motion Estimation for Nonoverlapping Multicamera Rigs: Linear Algebraic and Linf Geometric Solutions," *IEEE Trans. in Pattern Analysis and Machine Intelligence*, vol. 32, no. 6, pp. 1044–1058, 2010.
- [25] D. Nister, "An efficient solution to the five-point relative pose problem," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 26, no. 6, pp. 756–770, June 2004.
- [26] R. Hartley and P. Sturm, "Triangulation," *Computer Vision and Image Understanding*, 1997.
- [27] P. Sturm and B. Triggs, "A factorization based algorithm for multi image projective structure and motion," in *European Conference in Computer Vision*, 1996.
- [28] N. Levi and M. Werman, "The viewing graph," in *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, vol. 1, June 2003, pp. 1–518–I–522 vol.1.
- [29] K. Josephson, M. Byrod, F. Kahl, and K. Åström, "Image-based localization using hybrid feature correspondences," in *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, 2007, pp. 1–8.
- [30] P. Sturm, S. Ramalingam, J.-P. Tardif, S. Gasparini, and J. a. Barreto, "Camera models and fundamental concepts used in geometric computer vision," *Found. Trends. Comput. Graph. Vis.*, vol. 6, no. 1–2, pp. 1–183, Jan. 2011. [Online]. Available: <http://dx.doi.org/10.1561/06000000023>
- [31] H. Pottmann and J. Wallner, *Computational line geometry*, 1st ed. Berlin: Springer Verlag, 2001.
- [32] J. P. Barreto, "General central projection systems: Modeling, calibration and visual servoing," Ph.D. dissertation, PhD Thesis, University of Coimbra, Coimbra, Portugal, 2004.
- [33] N. Snavely, S. M. Seitz, and R. Szeliski, "Photo tourism: exploring photo collections in 3d," in *ACM transactions on graphics (TOG)*, vol. 25, no. 3. ACM, 2006, pp. 835–846.
- [34] J. Allard, J.-S. Franco, C. Merrier, E. Boyer, and B. Raffin, "The grimage platform: A mixed reality environment for interactions," in *Computer Vision Systems, 2006 ICVS'06. IEEE International Conference on*. IEEE, 2006, pp. 46–46.
- [35] "4d repository," <http://4drepository.inrialpes.fr/pages/home>.
- [36] C. Wu, "Towards linear-time incremental structure from motion," in *2013 International Conference on 3D Vision-3DV 2013*. IEEE, 2013, pp. 127–134.
- [37] "Samantha," <http://www.diegm.uniud.it/fusiello/demo/samantha/>.



Francisco Vasconcelos received the PhD degree from the University of Coimbra, Portugal, in 2016. He is a Postdoctoral Associate Researcher at the University College London, United Kingdom, as part of the Centre for Medical Image Computing and the Surgical Robot Vision group. His current research interests focus on geometry problems in computer vision, medical imaging, and robotics.



João P. Barreto received the Ph.D. degrees from the University of Coimbra, Coimbra, Portugal, in 2004. From 2003 to 2004, he was a Postdoctoral Researcher with the University of Pennsylvania, Philadelphia. Since 2004, he has been an Assistant Professor with the University of Coimbra, where he is also a Senior Researcher with the Institute for Systems and Robotics. He is a cofounder of the Perceive3D Company in the domain of image guided surgery. His current research

interests include different topics in computer vision, with a special emphasis in geometry problems and applications in robotics and medicine.



Edmond Boyer received the PhD degree from the Institut National Polytechnique de Lorraine, France, in 1996. He is an associate professor at Grenoble Universities, France. His fields of competence cover computer vision, computational geometry, and virtual reality. He is a cofounder of the 4D View Solution Company in the domain of spatiotemporal modeling. His current research interests are in 3D dynamic modeling from images and videos, motion capture and recognition from

videos, and immersive and interactive environments.