

Trifocal Relative Pose from Lines at Points and its Efficient Solution

Ricardo Fabbri

Rio de Janeiro State University

rfabbri@iprj.uerj.br

Timothy Duff

GA Tech

Hongyi Fan

Brown University

Margaret Regan

University of Notre Dame

David da Costa de Pinho

Northern Rio de Janeiro State University

UENF – Brazil

Elias Tsigaridas

INRIA Paris

Charles Wampler

General Motors R&D

Jonathan Hauenstein

University of Notre Dame

Peter J. Giblin

University of Liverpool

Benjamin Kimia

Brown University

Anton Leykin

GA Tech

Tomas Pajdla

CTU

Abstract

We present a new minimal problem for relative pose estimation mixing point features with lines incident at points observed in three views and its efficient homotopy continuation solver. We demonstrate the generality of the approach by analyzing and solving an additional problem with mixed point and line correspondences in three views. The minimal problems include correspondences of (i) three points and one line and (ii) three points and two lines through two of the points which is reported and analyzed here for the first time. These are difficult to solve, as they have 216 and – as shown here – 312 solutions, but cover important practical situations when line and point features appear together, e.g., in urban scenes or when observing curves. We demonstrate that even such difficult problems can be solved robustly using a suitable homotopy continuation technique and we provide an implementation optimized for minimal problems that can be integrated into engineering applications. Our simulated and real experiments demonstrate our solvers in the camera geometry computation task in structure from motion. We show that new solvers allow for reconstructing challenging scenes where the standard two-view initialization of structure from motion fails.

1. Introduction

Three-dimensional computer vision has made a wider impact [4], in part by relying on point-based structure from motion (SfM) [1, 63]. Matching point features across views leads to successful pose estimation and unorganized 3D point cloud reconstructions [49, 18]. Even production-quality SfM technology nevertheless fails [4] when the images contain (i) large homogeneous areas with few or no features; (ii) repeated textures, like brick walls, giving rise to a large number of ambiguously correlated features; (iii)

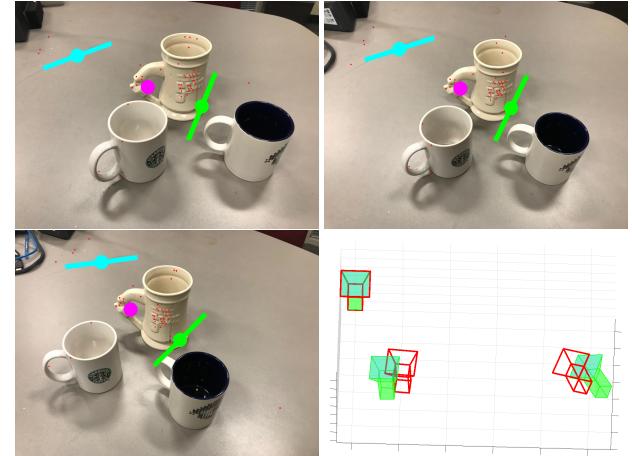


Figure 1. Three mug images illustrate deficiencies of the traditional two-view approach to bootstrapping SfM: there is not enough features detected and thus a SOTA SfM pipeline COLMAP [63] fails to reconstruct the relative pose of the cameras. In contrast, the truly trinocular method proposed here used the two triplets of point-tangents and one triplet of points (highlighted among red features) to reconstruct the pose of the cameras as shown. The schematic shows the matching of two triplets of points with attached lines in green and cyan and one point triplet without lines in pink. Red cameras were computed by our solver, and green cameras are the ground truth.

blurred areas, arising from fast moving cameras or objects; (iv) large scale changes where the feature overlap is not sufficiently significant; (v) Multiple and independently moving objects each of which do not have a sufficient number of features.

We track the failure cases to two key observations. First, multiview applications rarely make use of the full information available in the image sequence. Most traditional multiview pipelines estimate the relative pose of two views, say with the 5-point algorithm [47] and then register new views using a P3P algorithm [64]. Camera estimation from

trifocal tensors is believed to augment two-view pose estimation [16], although this is questioned in practice [28]. The calibrated trinocular relative pose estimation from four points, 3v4p, is known to be difficult to solve [48, 57, 58], partly because it is not a minimal problem since it is over-constrained. The first working solver [48] is effectively determining relative pose between two cameras in the form of a curve of degree ten of possible epipoles and using a third view to select the one that minimizes reprojection errors. In this sense, trinocular pose estimation has not truly been tackled as a minimal problem.

The second key observation is that low number of point features in images may often be supported by lines and curves. However, the use of points on curves to establish correspondence faces its own challenges. They are only transversally localized, leaving thus a dimension of ambiguity in determining curve correspondence. Despite this, curve points offer additional useful constraint, *i.e.* the orientation of their tangent. Thus, at corners, junctions, and other special points on curves, e.g., satisfying certain appearance conditions (maximizing the cornerness or Laplacian of Gaussian along the curve), enough points are both spatially localized, and orientation also available for an additional constraint. Of course, the availability of orientation is not exclusive to tangents on a curve; for example, we show how the SIFT dominant direction can be used effectively as an orientation attached to a point. We show that the introduction of “orientation attached to a point” can solve for estimation with fewer point matches, from 4 to 3, which is critical in images experiencing a feature drought, Figure 1, as well as to enhance the robustness and speed in RANSAC.

The two above observations motivate exploring trinocular pose estimation from the perspective of triplet point correspondences where the points may also be endowed with orientation. We demonstrate that only three points are needed when matched across all three views: Two of these triplets need to have attached orientation; the third does not; see the schematic in Figure 1.

Three types of constraints arise in matching points with attached orientation. First, the point location correspondence, *i.e.*, the epipolar constraint, provides an equation for each pair of views, or six equations in all. The fact that a pair from view 1 to view 2 and a pair from 2 to view 3 form a triplet provides another equation which essentially constrains the independent pairs of scale ambiguities to a single one. This provides another three equations. Finally, for each triplet of points with attached orientations, the orientation of the first two views predicts an orientation for the third, providing an additional constraint for each triplet with orientation. This provides two more equations, for a total of 11 equations in 11 unknowns.

These equations are polynomial with such complexity

that is not trivial to solve efficiently. This motivates using techniques from numerical algebraic geometry [7, 11, 40] to (i) probe whether the system is over or under constrained or otherwise minimal; (ii) understand the range of the number of solutions and a tight upper bound on it; (iii) develop efficient and practically relevant methods for finding solutions which are real and represent camera configurations. This paper answers all three points: the problem posed is minimal, it has up to 312 solutions of which 2-3 end up becoming relevant to camera configurations, and the paper develops a practical and relatively fast method (currently under 2 seconds but promises to be sub-second with some optimization) for solving the system; these are the key contribution of this paper. As a bonus, a similar trifocal problem with three points and a free line is analyzed to demonstrate generality of this approach.

Experiments are conducted on synthetic data to understand how the approach behaves under (i) veridical and accurate correspondences, (ii) veridical but noisy correspondences, and (iii) veridical noisy correspondences embedded among outliers. These experiments demonstrate that the system is robust and stable under spatial and orientation noise and under a significant level of outliers. For experiments on real data, we use SIFT keypoints endowed with SIFT orientation. The approach applies RANSAC to triplets of features, which are essentially pairs that are cycle consistent across three views, and solves the system of polynomial equations using an efficient implementation of homotopy continuation. The results are validated by measuring inliers. We have found that our approach is successful in all cases where the traditional SfM pipeline succeeds but more importantly it succeeds in many other cases too, on the EPFL [67] and Amsterdam Teahouse datasets [68], Figures 1 and 9. For additional details, we refer the reader to the supplementary material.

1.1. Literature Review

Trifocal Geometry Calibrated trifocal geometry estimation is a hard problem [57, 58, 48, 60]. There are no publicly available solvers we are aware of. The state of the art solver [48], based on four corresponding points (3v4p), has not yet found many practical applications [35].

For the uncalibrated case, 6 points are needed [21], and Larsson et al. recently solved the longstanding trifocal minimal problem of using 9 lines [36]. The case of mixed points and lines is less common [51, 51], but has seen a growing interest in related problems [69, 56]. The calibrated cases beyond 3v4p are largely unsolved, spurring more sophisticated theoretical work [30, 39, 42, 43, 2, 50, 3]. Kileel [30] studied many minimal problems in this setting, such as the Cleveland problem solved in the present paper, and reported studies using homotopy continuation. Kileel also stated that the full set of ideal generators is currently unknown, *i.e.*, a

given set of polynomial equations provably necessary and sufficient to describe calibrated trifocal geometry.

Seminal works used curves and edges in three views to transfer differential geometry for matching [5, 59], and for pose and trifocal tensor estimation [10, 62]. Point-tangents can be framed as *quivers* (1-quivers), or feature points with attributed directions (*e.g.*, corners), proposed in the context of uncalibrated trifocal geometry but de-emphasizing the connection to tangents to general curves [26, 71]. We note that point-tangent fields may also be framed as vector fields, so related technology may apply to surface-induced correspondence data [13]. In the calibrated setting, point-tangents were first used for absolute pose estimation by Fabbri *et al.* [14], using only two points, later relaxed for unknown focal length [34]. The trifocal problem with three point-tangents as a local version of trifocal pose for global curves was first formulated by Fabbri [13], for which we here present a minimal version codenamed Chicago.

Homotopy Continuation The basic theory of Polynomial Homotopy Continuation (HC) [7, 45, 65] was developed in 1976, and guarantees algorithms that are *globally* convergent with probability one from given start solutions. A number of general-purpose HC software have considerably evolved over the past decade [6, 9, 40, 70]. The computer vision community has used HC most notably in the nineties for 3D vision of curves and surfaces for tasks such as computing 3D line drawings from surface intersections, finding the stable singularities of a 3D line drawing under projections, computing occluding contours, stable poses, hidden line removal by continuation from singularities, aspect graphs, self-calibration, pose estimation [33, 53, 32, 53, 33, 52, 24, 8, 23, 44, 41, 17, 22, 55], as well as for MRFs [46, 8], and in more recent work [20, 12, 61]. An implementation of the early continuation solver of Kriegman and Ponce [32] by Pollefeys is still widely available for low degree systems [54].

As an early example [22], HC was used to find an early bound of 600 solutions to trifocal pose with 6 lines. In the vision community HC is mostly used as an offline tool to carry out studies of a problem before crafting a symbolic solver. Kasten *et al.* [29] recently compare a general purpose HC solver [70] against their symbolic solver. However, their problem is one order of magnitude lower degree than the ones presented here, and the HC technique chosen for our solver [11] is more specific than their use of polyhedral homotopy, in the sense that fewer paths are tracked (*c.f.* the start system hierarchy in [65]).

2. Two Trifocal Minimal Problems

2.1. Basic Equations

Our notation follows [19] with explicit projective scales. A more elaborate notation [10, 14] can be used to express the equations in terms of tangents to curves.

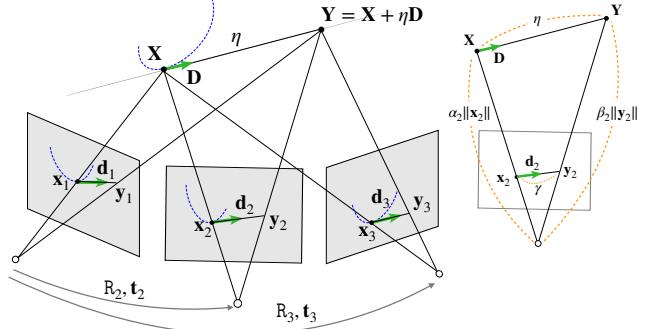


Figure 2. Notation for the trifocal pose problems.

Notation Let \mathbf{X} and \mathbf{Y} denote inhomogeneous coordinates of 3D points and $\mathbf{x}_{v,p} \in \mathbb{P}^2$ denote homogeneous coordinates of image points. Subscript v numbers views and p numbers the points. If only a single subscript is used, it indexes views. Symbols R_i, t_i denote the rotation and translation transforming coordinates from camera 1 to camera i , \mathbf{d} is an image line direction or curve tangent in homogeneous coordinates, and \mathbf{D} is the 3D line direction or space curve tangent in inhomogeneous world coordinates. Symbols α, β denote the depth of \mathbf{X}, \mathbf{Y} , respectively, and η is the displacement along \mathbf{D} corresponding to the displacement γ_i along \mathbf{d} .

We next formulate two minimal problems for points and lines in three views and derive their general equations before turning to specific formulations. We first state a new minimal problem codenamed ‘Chicago’, followed by an important similar problem, ‘Cleveland’.

Definition 1 (Chicago trifocal problem). Given three points $\mathbf{x}_{1v}, \mathbf{x}_{2v}, \mathbf{x}_{3v}$ and two lines ℓ_{1v}, ℓ_{1v} in views $v = 1, 2, 3$, such that the ℓ_{iv} meet \mathbf{x}_{iv} , $i = 1, 2$, $v = 1, 2, 3$, compute R_2, R_3, t_2, t_3 .

Definition 2 (Cleveland trifocal problem). Given three points $\mathbf{x}_{1v}, \mathbf{x}_{2v}, \mathbf{x}_{3v}$ in views $v = 1, 2, 3$, and given one line ℓ_{1v} in each image, compute R_2, R_3, t_2, t_3 .

To setup equations, we start with image projections of points $\alpha_1 \mathbf{x}_1 = \mathbf{X}, \alpha_2 \mathbf{x}_2 = R_2 \mathbf{X} + t_2, \alpha_3 \mathbf{x}_3 = R_3 \mathbf{X} + t_3$ and eliminate \mathbf{X} to get

$$\alpha_v \mathbf{x}_v = R_v \alpha_1 \mathbf{x}_1 + \mathbf{t}_v, \quad v = 1, 2 \quad (1)$$

Lines in space through \mathbf{X} are modeled by their points $\mathbf{Y} = \mathbf{X} + \eta \mathbf{D}$ in direction \mathbf{D} from \mathbf{X} . Points \mathbf{Y} are projected to

images as $\beta_1 \mathbf{y}_1 = \mathbf{X} + \eta \mathbf{D}$, $\beta_2 \mathbf{y}_2 = \mathbf{R}_2(\mathbf{X} + \eta \mathbf{D}) + \mathbf{t}_2$, $\beta_3 \mathbf{y}_3 = \mathbf{R}_3(\mathbf{X} + \eta \mathbf{D}) + \mathbf{t}_3$. Eliminating \mathbf{X} gives

$$\begin{aligned}\beta_1 \mathbf{y}_1 &= \alpha_1 \mathbf{x}_1 + \eta \mathbf{D} \\ \beta_2 \mathbf{y}_2 &= \alpha_2 \mathbf{x}_2 + \eta \mathbf{R}_2 \mathbf{D} \\ \beta_3 \mathbf{y}_3 &= \alpha_3 \mathbf{x}_3 + \eta \mathbf{R}_3 \mathbf{D}\end{aligned}\quad (2)$$

The directions \mathbf{d}_i of lines in images, which are obtained as the projection of \mathbf{Y} minus that of \mathbf{X} , i.e.

$$\beta_i \gamma_i \mathbf{d}_i = \mathbf{y}_i - \mathbf{x}_i = \alpha_i \mathbf{x}_i + \eta \mathbf{D} - \mathbf{x}_i, \quad (3)$$

are substituted to (2). After eliminating \mathbf{D} we get

$$(\beta_v - \alpha_v) \mathbf{x}_v + \beta_v \gamma_v \mathbf{d}_v = \mathbf{R}_v ((\beta_1 - \alpha_1) \mathbf{x}_1 + \beta_1 \gamma_1 \mathbf{d}_1), \quad (4)$$

For $v = 1, 2$. To simplify notation further, we change variables as $\epsilon_i = \beta_i - \alpha_i$, $\mu_i = \beta_i \gamma_i$ and get

$$\epsilon_v \mathbf{x}_v + \mu_v \mathbf{d}_v = \mathbf{R}_v (\epsilon_1 \mathbf{x}_1 + \mu_1 \mathbf{d}_1), \quad (5)$$

for $v = 1, 2$. For Chicago, we have three times the point equations (1) and two times the tangent equations (5). There are 12 unknowns $\mathbf{R}_2, \mathbf{t}_2, \mathbf{R}_3, \mathbf{t}_3$, and 24 unknowns $\alpha_{pv}, \epsilon_{pv}, \mu_{pv}$.

For Cleveland we need to represent a free 3D line L in space. We write a general point of L as $\mathbf{P} + \lambda \mathbf{V}$, with a point \mathbf{P} on L , the direction \mathbf{V} of L and real λ . Considering a triplet of corresponding lines represented by their homogeneous coordinates ℓ_v , the homogeneous coordinates of the back-projected planes are obtained as $\pi_v = [\mathbf{R}_v \mid \mathbf{t}_v]^T \ell_v$. Now, all π_v have to contain \mathbf{P} and \mathbf{V} and thus

$$\text{rank} [[I \mid 0]^T \ell_v \mid [\mathbf{R}_2 \mid \mathbf{t}_2]^T \ell_v \mid [\mathbf{R}_3 \mid \mathbf{t}_3]^T \ell_v] < 3 \quad (6)$$

Equations 6 and 1 are basic equations of Cleveland.

Many ways how to proceed by elimination from basic equations of the problems are possible. A particular formulation based on vanishing minors for both Chicago and Cleveland, which produced our first working solver to Chicago, is described in 3.1.

2.2. Problem Analysis

A general camera pose problem is defined by a list of labeled features in each image, which are in correspondence. The image coordinates of each feature are given, and we are to determine the relative poses of the cameras. The concatenated list of all the features' coordinates from all cameras is a point in the image space Y , while the concatenated list of the features' locations in the world frame or camera 1 is a point in the world feature space W . Unless the scale of some feature is given, the scale of the relative translations is indeterminate, so relative translations are treated as a projective space. For N cameras, the combined

poses of cameras $2, \dots, N$ relative to camera 1 are a point in $SE(3)^{N-1}$. Let the pose space be X , the projectivized version of $SE(3)^{N-1}$, and so $\dim X = 6N - 7$. Given the 3D features and the camera poses, we can compute the image coordinates of the features, so we have a viewing map $V: W \times X \rightarrow Y$. A camera pose problem is: given $y \in Y$, find $(w, x) \in W \times X$ such that $V(w, x) = y$. The projection $\pi: (w, x) \mapsto x$ is the set of relative poses we seek.

Definition 3. A camera pose problem is minimal if $V: W \times X \rightarrow Y$ is invertible and nonsingular at a generic $y \in Y$.

A necessary condition for a map to be invertible and nonsingular is that the dimensions of its domain and range must be equal. Let us consider three kinds of features: a point, a point on a line (equivalently a point with tangent direction), and a free line (a line with no distinguished point on it). For each feature, say F , let C_F be the number of cameras that see it. The contributions to $\dim W$ and $\dim Y$ of each kind of feature are in the table below, where a point with a tangent counts as one point and one tangent. Thus, a point feature has several tangents if several lines intersect at it (sometimes called quiver).

Feature	$\dim W$	$\dim Y$
Point, P	3	$2 \cdot C_P$
Tangent, T	2	$1 \cdot C_T$
Free Line, L	4	$2 \cdot C_L$

Accordingly, summing the contributions to $\dim Y - \dim W$ for all the features, we have the following result.

Theorem 2.1. Let $\langle x \rangle \doteq \max(0, x)$. A necessary condition for a N -camera pose problem to be minimal is

$$\sum_P \langle 2C_P - 3 \rangle + \sum_T \langle C_T - 2 \rangle + \sum_L \langle 2C_L - 4 \rangle = 6N - 7. \quad (7)$$

For trifocal problems where all cameras see all features, i.e., $C_P = C_T = C_L = 3$, a pose problem with 3 feature points and 2 tangents meets condition (7). A pose problem with 3 feature points and 1 free line also meets the condition. Adding any new features to these problems will make them overconstrained, having $\dim Y > \dim W \times X$.

To demonstrate sufficiency, it enough to find $(w, x) \in W \times X$ where the Jacobian of $V(w, x)$ is full rank. Choosing a random point (w, x) and testing the Jacobian rank serves to establish nonsingularity with probability one. Such a test computed in floating point arithmetic is highly indicative but not rigorous unless one bounds floating-point error, which can be done using interval arithmetic, or exact arithmetic. A singular value decomposition of the Jacobian computed in floating point that shows that the Jacobian has a smallest singular value far from zero, can be taken as a numerical demonstration that the problem is minimal. Similarly, a careful calculation using techniques from numerical algebraic geometry can compute a full solution list in \mathbb{C} for

a randomly selected example and thereby produce a numerical demonstration of the algebraic degree of the problem. Using such techniques, we make the following claims with the caveat that they have been demonstrated numerically, not proven rigorously.

Theorem 2.2 (Numerical). *The Chicago trifocal problem is minimal with algebraic degree 312, and the Cleveland problem is minimal with algebraic degree 216.*

Proof. The previous paragraphs explain the numerical arguments involved, but the definite proof by computer involves symbolically computing the Gröbner basis over \mathbb{Q} , with special provisions, as discussed in supplementary material. \square

While this result is in agreement with degree counts for Cleveland in [30], the analysis of Chicago is novel as this problem is presented in this paper for the first time. See the supplementary material.

3. Homotopy Continuation Solver

In this section we describe our homotopy continuation solvers. In subsection 3.1 we reformulate the trifocal pose estimation problems as parametric polynomial systems in unknowns R_2, R_3, t_2, t_3 using the main specific equations that so far have produced our best results, while other formulations are discussed in supplementary material. We attribute relatively good run times to two factors. First, we use coefficient-parameter homotopy, outlined in 3.2, which naturally exploits the algebraic degree of the problem. Already with general-purpose software [6, 40], parameter homotopies are observed to solve the problems in a relatively efficient manner. Secondly, we optimize various aspects of the homotopy continuation routine, such as polynomial evaluation and numerical linear algebra. In subsection 3.3, we describe our optimized implementation in C++ which was used to do the computations.

3.1. Equations based on minors

One way of building a parametric homotopy continuation solver is to formulate the problems as follows. An instance of Chicago may be described by 5 *visible lines* in each view. We represent each line by its defining equation in homogeneous coordinates, *i.e.* as $\ell_{i,1}, \dots, \ell_{i,5} \in \mathbb{C}^{3 \times 1}$ for each $i \in \{1, 2, 3\}$. With the convention that the first three lines pass through the three pairs of points in each view and that the last two pass through associated point-tangent pairs, let

$$L_j = [[I | 0]^T \ell_{1,j} \quad [R_2 | t_2]^T \ell_{2,j} \quad [R_3 | t_3]^T \ell_{3,j}] \quad (8)$$

for each $j \in \{1, \dots, 5\}$. We enforce *line correspondences* by setting all 3×3 minors of each L_j equal to zero. Certain *common point constraints* must also be satisfied, *i.e.*, that the

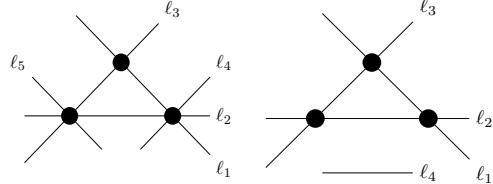


Figure 3. Visible line diagrams for Chicago and Cleveland.

4×4 minors of matrices $[L_1 | L_2 | L_4]$, $[L_1 | L_3 | L_5]$, and $[L_1 | L_2]$ all vanish.

We may describe the Cleveland problem with similar equations. For this problem, we are given lines $\ell_{i,1}, \dots, \ell_{i,4}$ for $i \in \{1, 2, 3\}$. We enforce line correspondences for matrices L_1, \dots, L_4 defined as in (8) and common point constraints by requiring that the 4×4 minors of $[L_1 | L_2]$, $[L_1 | L_3]$, and $[L_2 | L_3]$ all vanish. The “visible lines” representation of both problems is depicted in Figure 3.1.

3.2. Algorithm

From the previous section, we may define a specific system of polynomials $F(\mathcal{R}; \mathcal{A})$ in the unknowns $\mathcal{R} = (R_2, t_2, R_3, t_3)$ parametrized by $\mathcal{A} = (\ell_{1,1}, \dots)$. Many representations for rotations were explored, but our main implementation employs quaternions. A fundamental technique for solving such systems, fully described in [65], is *coefficient-parameter homotopy*. Algorithm 1 summarizes homotopy continuation from a known set of solutions for given parameter values to compute a set of solutions for the desired parameter values. It assumes that solutions for some starting parameters \mathcal{A}^* have already been computed via some offline, *ab initio* phase. For our problems of interest, the number of start solutions is precisely the algebraic degree of the problem.

Several techniques exist for the *ab initio solve*. For example, one can use standard homotopy continuation to solve the system $F(\mathcal{R}; \mathcal{A}^*) = 0$, where \mathcal{A}^* are randomly generated start parameters [65, 7]. This method may be enhanced by exploiting additional structure in the equations or using regeneration. Another technique based on monodromy, described in [11], was used to obtain a set of starting solutions and parameters for the solver described in Section 3.3.

3.3. Implementation

We provide an optimized C++ package called MINUS – MInimal problem NUmerical Solver <http://github.com/rfabbri/minus>. This is continuation code specialized for minimal problems, templated in C++, so that efficient specialization for different problems and different formulations are made possible. The most reliable and high-quality solver according to our experiments uses a 14×14 minors-based formulation. Although other formulations have demonstrated further potential for speedup by orders

Algorithm 1: Homotopy continuation solution tracker

input: Polynomial system $F(\mathcal{R}; \mathcal{A})$, where $\mathcal{R} = (\mathbf{R}_2, \mathbf{t}_2, \mathbf{R}_3, \mathbf{t}_3)$, and \mathcal{A} parametrizes the data;
 Start parameters \mathcal{A}^* ; start solutions \mathcal{R}^* where $F(\mathcal{R}^*; \mathcal{A}^*) = 0$; Target parameters $\widehat{\mathcal{A}}$
output: Set of target solutions $\widehat{\mathcal{R}}$ where $F(\widehat{\mathcal{R}}; \widehat{\mathcal{A}}) = 0$

Setup homotopy $H(\mathcal{R}; s) = F(\mathcal{R}; (1 - s)\mathcal{A}^* + s\widehat{\mathcal{A}})$.
for each start solution **do**
 $s \leftarrow \emptyset$
 while $s < 1$ **do**
 Select step size $\Delta s \in (0, 1 - s]$.
 Predict: Runge-Kutta Step from s to $s + \Delta s$ such that $dH/ds = 0$.
 Correct: Newton step st. $H(\mathcal{R}; s + \Delta s) = 0$.
 $s \leftarrow s + \Delta s$
return Computed solutions $\widehat{\mathcal{R}}$ where $H(\widehat{\mathcal{R}}, 1) = 0$.

of magnitude, there may be reliability tradeoffs (*c.f.* supplementary material).

4. Experiments

We first study the quality of our solver in synthetic experiments. Then, we demonstrate its performance on challenging real data. Due to space constraints, we present results for the Chicago problem, which is more challenging than Cleveland. See the supplementary material for experiments on Cleveland.

Synthetic experiments We show the performance of our solvers by starting with perfect synthetic data [15], consisting of 3D curves in a $4 \times 4 \times 4\text{cm}^3$ volume projected to 100 cameras Fig. 4, and sampling them to get 5117 potential data points/tangents that are projections of the same 3D analytic points and tangents [15], and then degrading them with noise and mismatches. Camera centers are randomly sampled around an average sphere around the scene along normally distributed radii of mean 1m and $\sigma = 10\text{mm}$, and rotations constructed via normally distributed look-at directions with mean along the sphere radius looking to the object, and $\sigma = 0.01$ rad such that the scene does not leave the viewport, followed by uniformly distributed roll. This sampling is filtered such that no two cameras are within 15° of each other.

Our first experiment studies the numerical stability of the solvers. The dataset provides true point correspondences, which inherit an orientation from the tangent to the analytic curve. For each sample set, three triplets of point correspondences are randomly selected with two endowed with the orientation of the tangent to the curve. The real solutions are selected from among the output, and only those that generate positive depth are retained. Finally, the un-

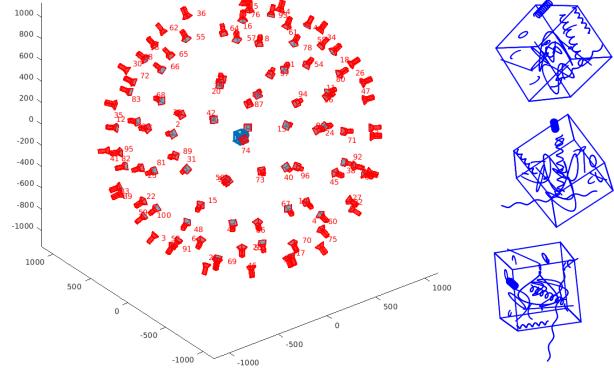


Figure 4. Sample views of our synthetic dataset. Real datasets have also been used in our experiments. (3D curves are from [14, 15]).

used tangent of the third triplet is used to verify the solution as it is an overconstrained problem. For each of the remaining solutions a pose is determined.

The error in pose estimation is compared with the ground truth, measured as the angular error between normalized translation vectors and the angular error between the quaternion representation of the rotation matrices. The entire process of generating the input to computing pose is repeated 1000 times and averaged. This experiment demonstrates that: (i) pose estimation errors are negligible, Fig. 5(a); (ii) the number of solutions is small – 35 real solutions on average which then get pruned down to around 7 on average by enforcing positive depth. Using the unused tangent of the third point as a verification reduces the number of physically realizable solutions to about 3 or 4, Fig. 5(b); (iii) The solver fails in about 1% of cases. These cases are detectable and while not a problem for RANSAC, the solver can be rerun for that solution path with higher accuracy or more parameters at a higher computational cost.

The second experiment shows that we can reliably and accurately determine cameras pose with correct but noisy correspondences. Using the same dataset and a subset of the selection of three triplets of points and tangents – 200 in total – zero-mean Gaussian noise was added both to the feature locations with σ corresponding to $\{0.25, 0.5, 0.75, 1.0\}$ pixels in image and to the orientation of the tangents with $\sigma \in \{0.25, 0.5, 0.75, 1.0\}$. These selected magnitude of localization errors reflect the expected localization error of point features and the orientation error corresponds to the state of the art orientation measurements [31]. A RANSAC scheme determines the feature set with pose generating the highest number of inliers. The experiments indicate that the resulting translation and rotation errors are reasonable. Figure 6(top) shows how changes in the magnitude of feature localization error affect pose in terms of translation errors and rotation errors. We use orientation perturbation of 0.1 rad to simulate the error in real feature orientation. Figure 6(bottom) shows how the magnitude of orientation error

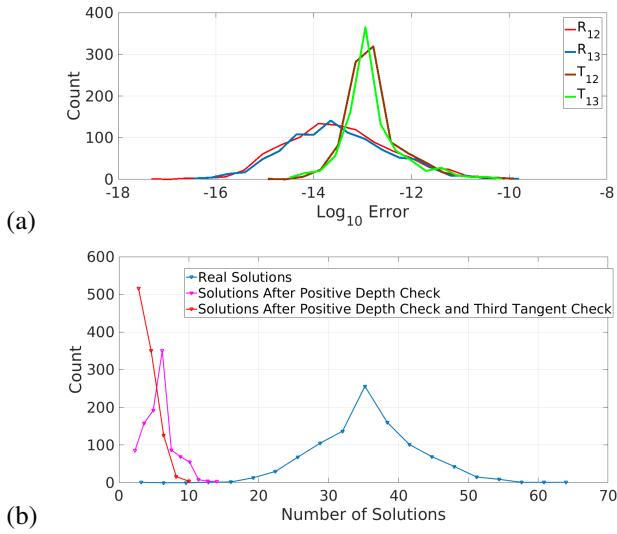


Figure 5. (a) Errors of computed parameters w.r.t the GT are small showing that the solver is numerically stable. (b) The histogram of the numbers of real solutions in different stages.

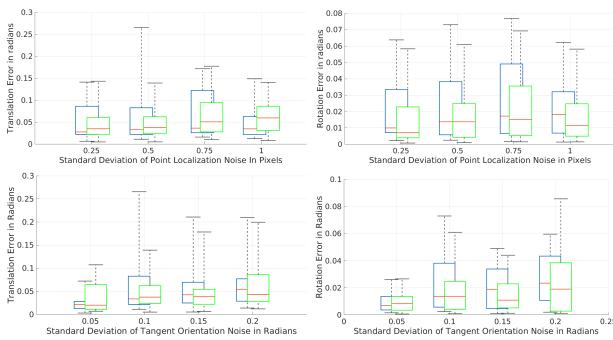


Figure 6. Distribution of trifocal pose error in the form of translational and rotational error is plotted against the level of feature localization noise and orientation noise. The green, resp., blue plots refer to the pose of the second resp., the third camera, relative to the first.

affects pose in terms of translation errors and rotation errors. A localization error of 0.5 pixel is used as orientation error is varied.

More meaningful, however, is the error measured in observation space, *i.e.*, the reprojection error: in each triplet of features, the first two features are used to predict the location of the third and the distance between the reprojected feature and the third perturbed feature is the reprojection error. This process is repeated 100 times to generate Figure 7.

The third experiment probes whether the system can reliably and accurately determine trifocal pose when veridical noisy correspondences are mixed with outliers. With an error of 0.25 pixels and 0.1 radians, 200 triplets of features were first generated and a percentage of these replaced with samples with random location and orientation. The ratio of outliers is 10%, 25% and 40%. The experiment was re-

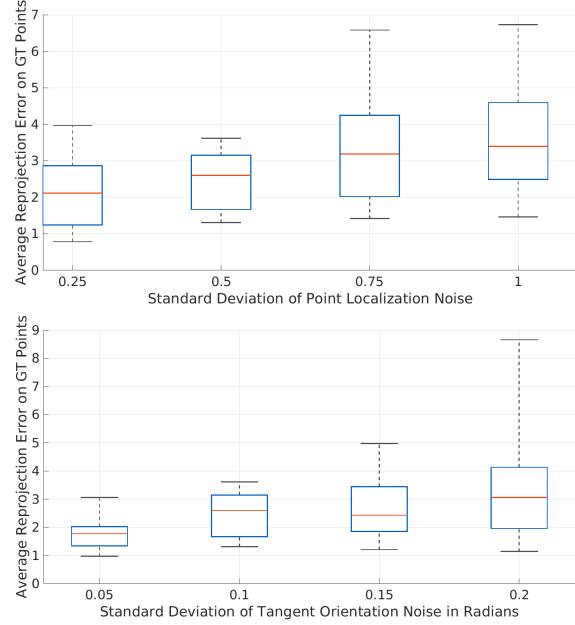


Figure 7. The distribution of reprojection error of feature location is plotted against of levels of feature localization error and feature orientation error.

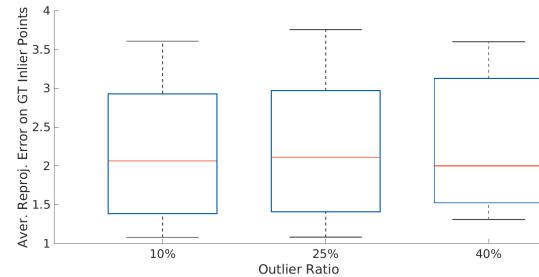


Figure 8. Average reprojection error on GT inlier points with different ratio of outliers.

peated 100 times. The resulting reprojection error is low and stable with the outlier ratio, Fig. 8.

Computational efficiency: Each step of minimal solve using our solver MINUS takes 1.9s in the worst case (about 660ms on average), corresponding to over 1 minute in our best prototypes using general purpose software [40, 6], both on an Intel core i7-7920HQ processor and four threads. More aggressive but potentially unsafe optimizations towards microseconds are feasible, but require assessing failure rate, as reported in the supplementary materials.

Real experiments: The use of attended lines in our approach requires harvesting points with attached tangents or orientations. In the case of isolated points, such as SIFT keypoints, the orientation of the SIFT descriptor allows a point to be endowed with an orientation. In the case of curves, the curve tangent provides a natural orientation for each point. However, while curve points show

superior transversal localization and superior orientation specification there is correspondence ambiguity along the curve. This can be resolved by employing corners and junctions [25] or special appearance-based keypoints found along a curve. One can also use the curve-to-curve correspondence ambiguity as part of a RANSAC procedure with some help from recent work [38]. These options are all viable. Since the focus of this paper is on the introduction of the approach, the solver, and a practical pipeline for trifocal pose estimation, we focus on the use of SIFT keypoints with SIFT orientations. We recognize that this is suboptimal, as the main drawback of feature-based relative pose estimation is in areas of low-number of features and repeated texture, so working with feature points inherits these difficulties. It would have been better to work with curves which are prominent and stable. Nevertheless we can use SIFT keypoints with attached SIFT orientations to illustrate that our method is at least as good as the traditional methods in all cases and in some cases solves the relative pose when the traditional scheme fails. It is worth emphasizing, however, that the potential of this scheme is to go beyond isolated features, a subject of future work.

Much like the standard pipeline, SIFT features are first extracted from all images. Pairwise features are found by rank-ordering measured similarities and making sure each feature's match in another image is not ambiguous and is above accepted similarity. Pairs of features from the first and second views are then grouped with the pairs of features from the second and third views into triplets. A cycle consistency check enforces that the triplets must also support a pair from the first and third views. Three feature triplets are then selected using RANSAC and together with their assigned SIFT orientation at two points used to estimate the relative pose of the three cameras.

Examples of this procedure are reported for triplets of images taken from the EPFL dense multi-view stereo test image dataset [66] in Figure 9, with ground-truth cameras shown in solid green and the cameras obtained with our method in red outlines. A qualitative visual comparison shows that our estimates are excellent. Quantitatively, our estimates have pose errors of 1.5×10^{-3} radians in translation and 3.24×10^{-4} radians in rotation. The average reprojection error is 0.310 pixels. These are comparable or better than the trifocal relative pose estimation methods reported in [27]. Our conclusion for this dataset is that our method is at least as good and often better than the traditional methods. See supplementary data for more examples and a substantiation of this claim.

The EPFL dataset, however, is texture-rich, typically yielding on the order of 1000 triplet features per triplet of images. As such it does not portray the typical problems faced in the really challenging situations when there are few features available or when there are repeated textures. The



Figure 9. Trifocal relative pose estimation of EPFL dataset. For each row, image triplets samples are shown. The estimation results are shown on the right. Ground truth poses are in solid green and estimated poses are in red. More examples in supplementary material.

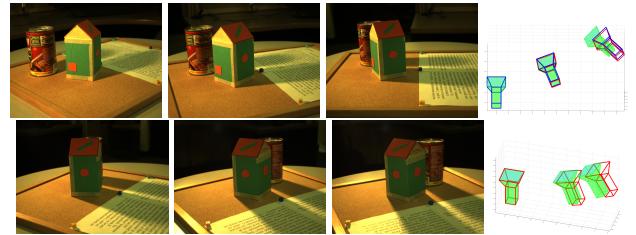


Figure 10. Two samples of trifocal relative camera pose estimation of Amsterdam Teahouse dataset. First line is a sample triplet of images that COLMAP is able to tackle. Second line is a sample triplet from the images that COLMAP reported it cannot find good matches. COLMAP results are in blue wireframes.

Amsterdam Teahouse Dataset [68], which also has ground-truth relative pose data, depicts scenes with fewer features. Figure 10 shows a triplet of images from this dataset where there is sufficient set of features (the soup can) to support a bifocal relative pose estimation followed by a P3P registration to a third view (using COLMAP [63].) However, when the number of features is reduced, as in Figure 10 where the number of features is much lower (soup can is invisible), COLMAP fails to find relative pose between pairs of these images. In contrast, our approach which relies on three and not five features is able to operate on this scene and recover the camera pose. Figure 1 shows another example. Further results are shown in supplementary material.

5. Conclusion

We presented a new calibrated trifocal minimal problem, an analysis demonstrating its number of solutions, and a practical solver by specializing numerical algebraic computation techniques. We show these techniques generalize to another difficult minimal problem with mixed points and lines. The proposed problem connects classical multi-view geometry of points and lines to that of points and tangents appearing when observing 3D curves extracted with tools of differential geometry [15, 13]. We believe that our ap-

proach to solving minimal problems may be useful for other difficult minimal problems. In the future, our “100 lines of custom-made solution tracking code” will be used to try to improve solvers of many other minimal problems which could not have been solved efficiently with Gröbner basis techniques [37].

References

- [1] S. Agarwal, N. Snavely, I. Simon, S. M. Seitz, and R. Szeliski. Building Rome in a day.
- [2] C. Aholt and L. Oeding. The ideal of the trifocal variety. *Math. Comp.*, 83, 2014.
- [3] A. Alzati and A. Tortora. A geometric approach to the trifocal tensor. *Journal of Mathematical Imaging and Vision*, 38(3):159–170, Nov 2010.
- [4] ARKit Team. Understanding ARKit tracking and detection. Apple, WWDC, 2018.
- [5] N. Ayache and L. Lustman. Fast and reliable passive binocular stereovision. In *1st International Conference on Computer Vision*, June 1987.
- [6] D. J. Bates, J. D. Hauenstein, A. J. Sommese, and C. W. Wampler. Bertini: Software for numerical algebraic geometry. Available at bertini.nd.edu.
- [7] D. J. Bates, J. D. Hauenstein, A. J. Sommese, and C. W. Wampler. *Numerically solving polynomial systems with Bertini*, volume 25 of *Software, Environments, and Tools*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2013.
- [8] A. M. Bruckstein, R. J. Holt, and A. N. Netravali. How to catch a crook. *J. Visual Communication and Image Representation*, 5(3):273–281, 1994.
- [9] T. Chen, T.-L. Lee, and T.-Y. Li. Hom4PS-3: A parallel numerical solver for systems of polynomial equations based on polyhedral homotopy continuation methods. In H. Hong and C. Yap, editors, *Mathematical Software – ICMS 2014*, pages 183–190, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg.
- [10] R. Cipolla and P. Giblin. *Visual Motion of Curves and Surfaces*. Cambridge University Press, 1999.
- [11] T. Duff, C. Hill, A. Jensen, K. Lee, A. Leykin, and J. Sommars. Solving polynomial systems via homotopy continuation and monodromy. *IMA Journal of Numerical Analysis*, 2018.
- [12] A. Ecker and A. D. Jepson. Polynomial shape from shading. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 145–152, June 2010.
- [13] R. Fabbri. *Multiview Differential Geometry in Application to Computer Vision*. Ph.D. dissertation, Division Of Engineering, Brown University, Providence, RI, 02912, July 2010.
- [14] R. Fabbri, P. J. Giblin, and B. B. Kimia. Camera pose estimation using first-order curve differential geometry. In *Proceedings of the IEEE European Conference in Computer Vision*, Lecture Notes in Computer Science. Springer, 2012.
- [15] R. Fabbri and B. B. Kimia. Multiview differential geometry of curves. *International Journal of Computer Vision*, 117:1–23, 2016.
- [16] O. Faugeras and Q.-T. Luong. *The Geometry of Multiple Images*. MIT Press, Cambridge, MA, USA, 2001.
- [17] O. D. Faugeras, Q. T. Luong, and S. J. Maybank. Camera self-calibration: Theory and experiments. In G. Sandini, editor, *Computer Vision — ECCV'92*, pages 321–334, Berlin, Heidelberg, 1992. Springer Berlin Heidelberg.
- [18] Y. Furukawa and J. Ponce. Accurate, dense, and robust multi-view stereopsis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32:1362–1376, 2010.
- [19] R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2nd edition, 2004.
- [20] J. D. Hauenstein and M. H. Regan. Adaptive strategies for solving parameterized systems using homotopy continuation. *Appl. Math. Comput.*, 332:19–34, 2018.
- [21] A. Heyden. Reconstruction from image sequences by means of relative depths. In *Proceedings of the Fifth International Conference on Computer Vision*, ICCV '95, pages 1058–, Washington, DC, USA, 1995. IEEE Computer Society.
- [22] R. J. Holt and A. N. Netravali. Motion and structure from line correspondences: Some further results. *International Journal of Imaging Systems and Technology*, 5(1):52–61, 1994.
- [23] R. J. Holt and A. N. Netravali. Number of solutions for motion and structure from multiple frame correspondence. *Int. J. Comput. Vision*, 23(1):5–15, May 1997.
- [24] R. J. Holt, A. N. Netravali, and T. S. Huang. Experience in using homotopy methods to solve motion estimation problems. volume 1251, 1990.
- [25] K. Huang, Y. Wang, Z. Zhou, T. Ding, S. Gao, and Y. Ma. Learning to parse wireframes in images of man-made environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 626–635, 2018.
- [26] B. Johansson, M. Oskarsson, and K. Astrom. Structure and motion estimation from complex features in three views. In *Proceedings of the Indian Conference on computer vision, graphics, and image processing*, 2002.
- [27] L. F. Julià and P. Monasse. A critical review of the trifocal tensor estimation. In *Pacific-Rim Symposium on Image and Video Technology*, pages 337–349. Springer, 2017.
- [28] L. Julià and P. Monasse. A critical review of the trifocal tensor estimation. In *The Eighth Pacific-Rim Symposium on Image and Video Technology – PSIVT'17*, Wuhan, China, 2017.
- [29] Y. Kasten, M. Galun, and R. Basri. Resultant based incremental recovery of camera pose from pairwise matches. *CoRR*, abs/1901.09364, 2019.
- [30] J. Kileel. Minimal problems for the calibrated trifocal variety. *SIAM Journal on Applied Algebra and Geometry*, 1(1):575–598, 2017.
- [31] B. B. Kimia, X. Li, Y. Guo, and A. Tamrakar. Differential geometry in edge detection: accurate estimation of position, orientation and curvature. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [32] D. J. Kriegman and J. Ponce. Curves and surfaces. chapter A New Curve Tracing Algorithm and Some Applications, pages 267–270. Academic Press Professional, Inc., San Diego, CA, USA, 1991.

- [33] D. J. Kriegman and J. Ponce. Geometric modeling for computer vision. volume 1610, 1992.
- [34] Y. Kuang and K. Åström. Pose estimation with unknown focal length using points, directions and lines. In *International Conference on Computer Vision*, pages 529–536. IEEE, 2013.
- [35] Y. Kuang, M. Oskarsson, and K. Åström. Revisiting trifocal tensor estimation using lines. In *Pattern Recognition (ICPR), 2014 22nd International Conference on*, pages 2419–2423. IEEE, 2014.
- [36] V. Larsson, K. Åström, and M. Oskarsson. Efficient solvers for minimal problems by syzygy-based reduction. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [37] V. Larsson, M. Oskarsson, K. Åström, A. Wallis, Z. Kukelova, and T. Pajdla. Beyond grobner bases: Basis selection for minimal solvers. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 3945–3954, 2018.
- [38] P. Lei, F. Li, and S. Todorovic. Joint spatio-temporal boundary detection and boundary flow prediction with a fully convolutional siamese network. *CVPR*, 2018.
- [39] S. Leonardos, R. Tron, and K. Daniilidis. A metric parametrization for trifocal tensors with non-colinear pinholes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 259–267, June 2015.
- [40] A. Leykin. Numerical algebraic geometry. *J. Softw. Alg. Geom.*, 3:5–10, 2011.
- [41] Q.-T. Luong. *Matrice Fondamentale et Calibration Visuelle sur l’Environnement-Vers une plus grande autonomie des systemes robotiques*. PhD thesis, Université de Paris-Sud, Centre d’Orsay, 1992.
- [42] E. Martyshev. On some properties of calibrated trifocal tensors. *Journal of Mathematical Imaging and Vision*, 58(2):321–332, 2017.
- [43] J. Mathews. Multi-focal tensors as invariant differential forms. *arXiv e-prints*, page arXiv:1610.04294, Oct 2016.
- [44] S. J. Maybank and O. D. Faugeras. A theory of self-calibration of a moving camera. *Int. J. Comput. Vision*, 8(2):123–151, 1992.
- [45] A. Morgan. *Solving polynomial systems using continuation for engineering and scientific problems*, volume 57 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2009. Reprint of the 1987 original.
- [46] P. K. Nanda, U. B. Desai, and P. Poonacha. A homotopy continuation method for parameter estimation in mrf models and image restoration. In *Proceedings of IEEE International Symposium on Circuits and Systems - ISCAS '94*, 1994.
- [47] D. Nister. An efficient solution to the five-point relative pose problem. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 26(6):756–770, 2004.
- [48] D. Nistér and F. Schaffalitzky. Four points in two or three calibrated views: Theory and practice. *Int. J. Comput. Vision*, 67(2):211–231, 2006.
- [49] D. Nistér, O. Naroditsky, and J. Bergen. Visual odometry. In *Computer Vision and Pattern Recognition (CVPR)*, pages 652–659, 2004.
- [50] L. Oeding. The quadrifocal variety. *arXiv e-prints*, 2015.
- [51] M. Oskarsson, A. Zisserman, and K. Astrom. Minimal projective reconstruction for combinations of points and lines in three views. *Image and Vision Computing*, 22(10):777 – 785, 2004. British Machine Vision Computing 2002.
- [52] S. Petitjean. Algebraic geometry and computer vision: Polynomial systems, real and complex roots. *Journal of Mathematical Imaging and Vision*, 10(3):191–220, May 1999.
- [53] S. Petitjean, J. Ponce, and D. J. Kriegman. Computing exact aspect graphs of curved objects: Algebraic surfaces. *International Journal of Computer Vision*, 9(3):231–255, Dec 1992.
- [54] M. Pollefeys. Vnl realnpoly: A solver to compute all the roots of a system of n polynomials in n variables through continuation. Available at https://github.com/vxl/vxl/blob/master/core/vnl/algo/vnl_rnlpoly_solve.h, 1997.
- [55] M. Pollefeys and L. Van Gool. Stratified self-calibration with the modulus constraint. *IEEE Trans. Pattern Anal. Mach. Intell.*, 21(8):707–724, Aug. 1999.
- [56] A. Qadir and J. Neubert. A line-point unified solution to relative camera pose estimation. *CoRR*, abs/1710.06495, 2017.
- [57] L. Quan, B. Triggs, and B. Mourrain. Some results on minimal euclidean reconstruction from four points. *J. Math. Imaging Vis.*, 24(3):341–348, 2006.
- [58] L. Quan, B. Triggs, B. Mourrain, and A. Ameller. Uniqueness of minimal Euclidean reconstruction from 4 points. Technical report, 2003. unpublished article.
- [59] L. Robert and O. D. Faugeras. Curve-based stereo: figural continuity and curvature. In *Proceedings of Computer Vision and Pattern Recognition*, pages 57–62, June 1991.
- [60] V. Rodehorst. Evaluation of the metric trifocal tensor for relative three-view orientation. In *International Conference on the Application of Computer Science and Mathematics in Architecture and Civil Engineering*, July 2015.
- [61] M. Salzmann. Continuous inference in graphical models with polynomial energies. In *CVPR*, pages 1744–1751. IEEE Computer Society, 2013.
- [62] C. Schmid and A. Zisserman. The geometry and matching of lines and curves over multiple views. *International Journal of Computer Vision*, 40(3):199–233, 2000.
- [63] J. L. Schönberger and J.-M. Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [64] N. Snavely, S. M. Seitz, and R. Szeliski. Modeling the world from internet photo collections. *International Journal of Computer Vision (IJCV)*, 80(2):189–210, 2008.
- [65] A. J. Sommese and C. W. Wampler, II. *The numerical solution of systems of polynomials: arising in engineering and science*. World Scientific Publishing Co. Pte. Ltd., Hackensack, NJ, 2005.
- [66] C. Strecha, W. von Hansen, L. J. V. Gool, P. Fua, and U. Thoennessen. On benchmarking camera calibration and multi-view stereo for high resolution imagery. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008), 24-26 June 2008, Anchorage, Alaska, USA*, 2008.

- [67] C. Strecha, W. von Hansen, L. Van Gool, P. Fua, and U. Thoennessen. On benchmarking camera calibration and multi-view stereo for high resolution imagery. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8, June 2008.
- [68] A. Usumezbas, R. Fabbri, and B. B. Kimia. From multiview image curves to 3D drawings. In *Proceedings of the European Conference in Computer Visiohn*, 2016.
- [69] A. Vakhitov, V. Lempitsky, and Y. Zheng. Stereo relative pose from line and point feature triplets. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [70] J. Verschelde. Algorithm 795: PHCpack: A general-purpose solver for polynomial systems by homotopy continuation. *ACM Trans. Math. Softw.*, 25(2):251–276, June 1999.
- [71] J. Zhao, L. Kneip, Y. He, and J. Ma. Minimal case relative pose computation using ray-point-ray features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2019.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. DMS-1439786 while most authors were in residence at Brown University’s Institute for Computational and Experimental Research in Mathematics – ICERM, in Providence, RI, during the Fall 2018 and Spring 2019 semesters.