

# Quaternion Based Camera Pose Estimation From Matched Feature Points

Kaveh Fathian, J. Pablo Ramirez-Paredes, Emily A. Doucette, J. Willard Curtis, Nicholas R. Gans.

**Abstract**—We present a novel solution to the camera pose estimation problem, where rotation and translation of a camera between two views are estimated from matched feature points in the images. The camera pose estimation problem is traditionally solved via algorithms that are based on the essential matrix or the Euclidean homography. With six or more feature points in general positions in the space, essential matrix based algorithms can recover a unique solution. However, such algorithms fail when points are on critical surfaces (e.g., coplanar points) and homography should be used instead. By formulating the problem in quaternions and decoupling the rotation and translation estimation, our proposed algorithm works for all point configurations. Using both simulated and real world images, we compare the estimation accuracy of our algorithm with some of the most commonly used algorithms. Our method is shown to be more robust to noise and outliers. For the benefit of community, we have made the implementation of our algorithm available online and free<sup>1</sup>.

## I. INTRODUCTION

Many applications in computer vision and robotics require measurements of the rotation and translation (i.e., *pose*) changes of an object as it moves through an environment. In photogrammetry, for example, by knowing the pose changes of the camera, 3D model of a scene can be constructed from a set of 2D images [1]. In robotics, pose estimated from images can be used for navigation [2], [3], or fused with other sensor measurements (e.g., IMU and GPS) to increase the reliability and accuracy [4]. Camera pose estimation has further applications in simultaneous localization and mapping (SLAM) [5], autonomous vehicles [6], and augmented reality [7].

Camera pose estimation techniques are often based on image features (e.g., edges, corners, etc., in an image) that can be detected and matched in two or more images [8], [9]. Figure 1 shows an example where feature points are detected and matched as indicated by yellow lines in two images. Many existing methods [10]–[13] use the coordinates of feature points on the image to construct the essential/fundamental matrix or the Euclidean homography

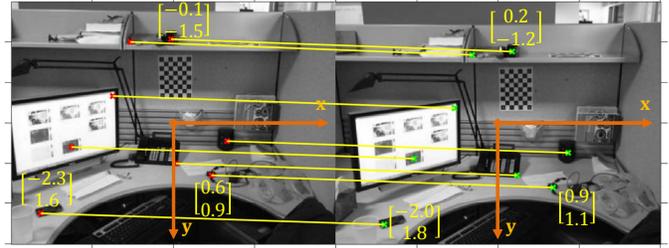


Fig. 1. Pairs of matched feature points across two images shown side by side (image courtesy of MathWorks).

matrix, from which relative rotation and translation of the camera can be recovered. Although these approaches are fast and easy to implement, they are subject to drawbacks: the essential matrix based algorithms fail when feature points are on critical surfaces (e.g., coplanar points), while homography based algorithms only work when points are coplanar.

In this work, we present a novel formulation of the camera pose estimation problem using quaternions and present a solution to estimate the pose under this formulation. Our approach, which we refer to as the Quaternion Estimation (QuEst) algorithm, does not use the homography or essential matrices, and decouples the estimation of rotation and translation. Consequently, common problems such as degeneracy for special 3D point configurations are avoided. We present two methods to recover the rotation from seven and six matched feature points. We then show how the unique correct solution can be detected from among the set of recovered solutions. The performance of QuEst is compared with algorithms that are based on the homography or essential matrix in the presence of noise in image point coordinates. The performance is further vetted by using real world image datasets that come with the ground truth camera pose information.

The main contributions and benefits of the proposed algorithm can be summarized as follows.

- As illustrated in Table I, unlike the homography or essential matrix based algorithms, QuEst provides an accurate pose estimate for both general point configurations and points that are on critical surfaces, e.g., coplanar points. To initialize the bundle adjustment in SLAM applications [5], often heuristic methods are used to detect the coplanarity of the points and select the appropriate algorithm correspondingly. QuEst can be used to initialize the bundle adjustment regardless of the feature point configuration in the space.
- By recovering the translation, QuEst simultaneously

<sup>\*</sup>This work was supported by the OSD sponsored Autonomy Research Pilot Initiative project entitled A Privileged Sensing Framework.

K. Fathian and N. R. Gans are with the Department of Electrical Engineering, University of Texas at Dallas, Richardson, TX, 75080 USA. E-mail: {kaveh.fathian, ngans}@utdallas.edu.

JP Ramirez is with the Research and Engineering Education Facility, University of Florida, Shalimar, FL, 32579, USA. E-mail: : jpramirez@ufl.edu.

E. A. Doucette and J. W. Curtis are with the Air Force Research Laboratory, Munitions Directorate, Eglin AFB, FL, 32542, USA. E-mail: {emily.doucette, jess.curtis}@us.af.mil.

<sup>1</sup>The Matlab implementation of QuEst is accessible at <https://goo.gl/QH5qhw>

TABLE I. QuEst compared with homography and essential matrix based algorithms.

Approach	Advantage	Works for general points	Works for coplanar points
Essential matrix (8-6 point)		✓	✗
Euclidean homography		✗	✓
QuEst (6 point)		✓	✓

recovers depths of the feature points. Therefore, a 3D model of the scene can be reconstructed. The recovered translation and depths share a common scale factor, hence, the magnitude of the recovered translation vector goes to zero as the camera translation between two views approaches zero. The translation recovered from the essential matrix always has unit norm, which is not desirable in applications such as visual servoing.

- Tests performed using both simulated and real world images show that the pose recovered from QuEst is more accurate and robust to noise and outliers compared to the existing algorithms.

The rest of the paper is organized as follows. We briefly review related approaches in Section II, before we formulate the pose estimation problem using quaternions in Section III. We present the QuEst algorithm in Section IV and evaluate its performance under noise in Section V. In Section IV, we further vet the performance of QuEst using the real world image datasets.

## II. RELATED WORK

Pose estimation using images can be traced back to at least 1913, when Kruppa [14] proved that two camera views of five 3D points could be used to estimate the translation and rotation separating two camera views. However, the lack of sufficient computational resources limited further development until the landmark introduction of the 8-point algorithm [10], [11]. The 8-point algorithm used the epipolar constraint, embodied in the essential matrix, to solve a set of linear equations that are generated from a set of 8 or more feature points.

It has been shown that to ensure a finite number of solutions for the pose estimation problem, a minimum of five feature points are required [11], [15]. However, in general this is not sufficient to recover the pose uniquely. Philip showed that when 6 or more general points are available, the pose estimation problem is linear and produces a unique solution [16].

What previous work have in common is the use of the essential matrix in the problem formulation and the assumption that points are in general positions. In the special case where feature points are coplanar, methods based on the homography matrix should be used instead to estimate the pose [12], [13]. However, for a general point configuration that does not contain at least four coplanar points, homography does not return the correct solution.

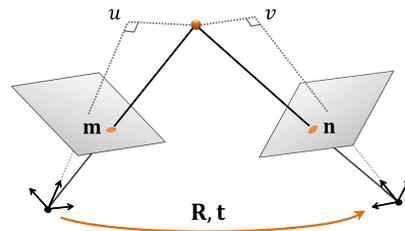


Fig. 2. Projection of a 3D point onto the image plane at two views.

We previously used quaternion formulation to solve the pose estimation via an optimization-based method [17]. The drawback of this initial work was its reliance on an accurate initial estimate for the pose. The method proposed here is based on formulating an eigenvalue problem, which does not require an initial guess for the solution.

## III. CAMERA POSE ESTIMATION

We introduce the notation and assumptions, followed by formulating the pose estimation problem in quaternions.

### A. Notation and Assumptions

Throughout the paper, we assume that the camera calibration matrix is known; this matrix can be easily found through the existing camera calibration routines [18].

Scalars are represented by lower case (e.g.,  $s$ ), vectors by lowercase and bold (e.g.,  $\mathbf{v}$ ), and matrices by upper case and bold letters (e.g.,  $\mathbf{M}$ ). All vectors are column vectors. The Moore-Penrose pseudo inverse of matrix  $\mathbf{M}$  is shown by  $\mathbf{M}^\dagger$ . Binomial coefficients are denoted by  $\binom{n}{k} := \frac{n!}{k!(n-k)!}$ . By norm of a vector, we imply the  $l^2$ -norm. Degree 4 *monomials* in variables  $w, x, y, z$  are single term polynomials  $w^a x^b y^c z^d$ , such that  $a + b + c + d = 4$ , for  $a, b, c, d \in \{0, 1, 2, 3, 4\}$ .

### B. Problem Formulation

Consider images of a scene taken by a camera at two views (e.g., Fig. 1). Let  $\mathbf{R} \in \text{SO}(3)$  and  $\mathbf{t} \in \mathbb{R}^3$  respectively represent the relative rotation and translation of the camera frame between the views. Assume that feature points are detected and matched, and their  $x$ - $y$  coordinates are read from the images. (In practice, the coordinates are in pixels, and should be mapped via the camera calibration matrix to Cartesian coordinates on the image plane.)

For each matched feature point the rigid motion constraint

$$u\mathbf{R}\mathbf{m} + \mathbf{t} = v\mathbf{n} \quad (1)$$

must hold, in which  $\mathbf{m}, \mathbf{n} \in \mathbb{R}^3$  are homogeneous coordinates (i.e., a 1 is appended to the  $x$ - $y$  coordinates) of the feature point in two images. Scalars  $u$  and  $v$  represent depths of the 3D point at each view, and as shown in Fig. 2, are the projections of the point onto the  $z$ -axis of the camera coordinate frame. Point coordinates  $\mathbf{m}$  and  $\mathbf{n}$  are known from the images, and the unknowns in (1) are  $u, v, \mathbf{R}$ , and  $\mathbf{t}$ , which need to be recovered.

We need to point out that in the pose estimation problem, translation and depths of the points can only be recovered up to a scale factor. This can be seen from (1), where any

constant multiplied into both hand sides can be absorbed by unknown variables  $u$ ,  $v$ , and  $\mathbf{t}$ .

**Example 1.** Consider pictures shown in Fig. 1, where feature points are matched and their coordinates on the image are determined. For the matched feature point with coordinates  $(-0.1, -1.5)$  in the left image and  $(0.2, -1.2)$  in the right image, from (1) we get

$$u_1 \mathbf{R} \begin{bmatrix} -0.1 \\ -1.5 \\ 1 \end{bmatrix} + \mathbf{t} = v_1 \begin{bmatrix} 0.2 \\ -1.2 \\ 1 \end{bmatrix}, \quad (2)$$

where scalar  $u_1$  and  $v_1$  are the depths of the 3D point at each view. Similarly, for two other matched pairs we can write

$$u_2 \mathbf{R} \begin{bmatrix} -2.3 \\ 1.6 \\ 1 \end{bmatrix} + \mathbf{t} = v_2 \begin{bmatrix} -2.0 \\ 1.8 \\ 1 \end{bmatrix}, \quad (3)$$

$$u_3 \mathbf{R} \begin{bmatrix} 0.6 \\ 0.9 \\ 1 \end{bmatrix} + \mathbf{t} = v_3 \begin{bmatrix} 0.9 \\ 1.1 \\ 1 \end{bmatrix}, \quad (4)$$

where subscripts are used to distinguish the depths of the points (scalars  $u$  and  $v$ ). Notice that rotation matrix  $\mathbf{R}$  and translation vector  $\mathbf{t}$  are the same in all equations.

Equation (1) uses the matrix representation of rotation, in which  $\mathbf{R}$  is a  $3 \times 3$  orthonormal matrix. That is,  $\mathbf{R}^\top = \mathbf{R}^{-1}$ , and  $\det(\mathbf{R}) = 1$ . Representing the rotation in the matrix form with orthonormality constraints makes the problem very nonlinear and challenging to solve. Instead, we use quaternions, which represent a rotation by four elements  $w, x, y, z \in \mathbb{R}$  such that  $w^2 + x^2 + y^2 + z^2 = 1$ . Although (1) can be formulated directly in quaternions, for simplicity and to avoid introducing the quaternion algebra we only mention what is essential to solve the problem here: if  $w, x, y, z \in \mathbb{R}$  are elements of a rotation quaternion, the associated rotation matrix is given by

$$\mathbf{R} = \begin{bmatrix} w^2 + x^2 - y^2 - z^2 & 2(xy - wz) & 2(xz + wy) \\ 2(xy + wz) & w^2 - x^2 + y^2 - z^2 & 2(yz - wx) \\ 2(xz - wy) & 2(yz + wx) & w^2 - x^2 - y^2 + z^2 \end{bmatrix}. \quad (5)$$

Quaternions provide a singularity free representation of rotation, and by restricting the first element to nonnegative numbers (i.e.,  $w \geq 0$ ), there is a one to one and onto correspondence between rotation matrices and quaternions.

To recover the pose, we first eliminate the unknowns  $u$ ,  $v$  and  $\mathbf{t}$ , and derive a system of equations in terms of the quaternion elements. From solving this system all rotation solution candidates are found. Subsequently, the translation and depths of the points are recovered. By taking (1) for two different feature points and subtracting the equations,  $\mathbf{t}$  can be eliminated. Subsequently,  $u$  and  $v$  can be eliminated from the resulting equations by noting that they form a null vector for the matrix consisting of the point coordinates and their rotations. The following example illustrates this procedure.

**Example 2.** Consider Example 1. By subtracting (3) and (4) from (2), respectively, and bringing terms to the left hand side we get

$$u_1 \mathbf{R} \begin{bmatrix} -0.1 \\ -1.5 \\ 1 \end{bmatrix} - v_1 \begin{bmatrix} 0.2 \\ -1.2 \\ 1 \end{bmatrix} - u_2 \mathbf{R} \begin{bmatrix} -2.3 \\ 1.6 \\ 1 \end{bmatrix} + v_2 \begin{bmatrix} -2.0 \\ 1.8 \\ 1 \end{bmatrix} = \mathbf{0}, \quad (6)$$

$$u_1 \mathbf{R} \begin{bmatrix} -0.1 \\ -1.5 \\ 1 \end{bmatrix} - v_1 \begin{bmatrix} 0.2 \\ -1.2 \\ 1 \end{bmatrix} - u_3 \mathbf{R} \begin{bmatrix} 0.6 \\ 0.9 \\ 1 \end{bmatrix} + v_3 \begin{bmatrix} 0.9 \\ 1.1 \\ 1 \end{bmatrix} = \mathbf{0}, \quad (7)$$

in which the unknown translation  $\mathbf{t}$  has been eliminated. We can represent (6) and (7) in the matrix-vector form

$$\underbrace{\begin{bmatrix} \mathbf{R} \begin{bmatrix} -0.1 \\ -1.5 \\ 1 \end{bmatrix} & \begin{bmatrix} -0.2 \\ 1.2 \\ 1 \end{bmatrix} & \mathbf{R} \begin{bmatrix} 2.3 \\ -1.6 \\ 1 \end{bmatrix} & \begin{bmatrix} -2.0 \\ 1.8 \\ 1 \end{bmatrix} & \mathbf{0} & \mathbf{0} \\ \mathbf{R} \begin{bmatrix} -2.0 \\ 1.8 \\ 1 \end{bmatrix} & \begin{bmatrix} -0.2 \\ 1.2 \\ 1 \end{bmatrix} & \mathbf{0} & \mathbf{0} & \mathbf{R} \begin{bmatrix} -0.6 \\ -0.9 \\ 1 \end{bmatrix} & \begin{bmatrix} 0.9 \\ 1.1 \\ 1 \end{bmatrix} \end{bmatrix}}_{\mathbf{M}} \begin{bmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \end{bmatrix} = \mathbf{0}, \quad (8)$$

which implies that matrix  $\mathbf{M} \in \mathbb{R}^{6 \times 6}$  has a null vector. Therefore, its determinant must be zero. Calculating determinant of  $\mathbf{M}$  with  $\mathbf{R}$  given in the parametric form (5) gives

$$\begin{aligned} & -0.1w^4 + 2.4w^3x + 35.6w^2x^2 - 11.4wx^3 + 0.3x^4 + \dots \\ & -19.8y^2z^2 - 168wz^3 + 9.4xz^3 - 17.8yz^3 + 0.3z^4 = 0. \end{aligned} \quad (9)$$

Equation (9) consists of all degree 4 monomials in  $w, x, y, z$  (i.e., terms such as  $w^4, w^3x, w^2x^2, \dots$ ), with coefficients that depend on the feature point coordinates. Note that since  $\mathbf{M}$  is a  $6 \times 6$  matrix, one may expect its determinant to have degree six monomials, however, due to special structure of  $\mathbf{M}$ , it is always possible to factor out  $w^2 + x^2 + y^2 + z^2$  from the determinant expression. Since  $w^2 + x^2 + y^2 + z^2 = 1$ , the degree four polynomial equation (9) follows.

What we showed in Example 2 was that three matched feature points generate a polynomial equation of the form (9). This equation is in terms of degree four monomials in  $w, x, y, z$ . Note that there are 35 such monomials, and their coefficients are in terms of the feature point coordinates. In practice, we do not need to calculate the determinant of  $\mathbf{M}$  to find these coefficients. By replacing the feature point coordinates with symbolic expressions the determinant can be computed symbolically, and explicit formulas for the coefficients can be derived. By substituting the numerical values of point coordinates in these formulas the coefficients are calculated directly. Due to the space limitation we do not give the explicit formulas here.

Since any three feature points give a polynomial equation of the form (9), from  $n$  points  $\binom{n}{3}$  equations can be generated. These equations can be stacked into a matrix-vector form, where the vector consists of the unknown monomial terms. The following example illustrates this point.

**Example 3.** Consider  $\binom{6}{3} = 20$  polynomial equations of the form (9), generated from 6 feature points. These polynomial equations can be represented in the matrix-vector form

$$\underbrace{\begin{bmatrix} -0.1 & 2.4 & 35.6 & \dots & 0.3 \\ -0.1 & 0.1 & -7.6 & \dots & -0.2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0.0 & -0.2 & 4.6 & \dots & 0.1 \end{bmatrix}}_{\mathbf{A}} \underbrace{\begin{bmatrix} w^4 \\ w^3x \\ w^2x^2 \\ \vdots \\ z^4 \end{bmatrix}}_{\mathbf{x}} = \mathbf{0} \quad (10)$$

where the coefficient matrix  $\mathbf{A} \in \mathbb{R}^{20 \times 35}$  depends on the feature point coordinates, and vector  $\mathbf{x} \in \mathbb{R}^{35}$  consists of all degree 4 monomials. Our goal is to find all  $w, x, y, z$ , for which (10) is satisfied.

The problem of recovering the rotation is henceforth equivalent to solving a system of equations of the form (10), where the goal is to find all  $\mathbf{x}$  for which

$$\mathbf{A}\mathbf{x} = \mathbf{0} \quad (11)$$

is satisfied. In (11), the coefficient matrix  $\mathbf{A}$  is known from the feature point coordinates, and vector  $\mathbf{x}$  is unknown with entries in degree four monomials of  $w, x, y, z$ .

#### IV. THE QUEST ALGORITHM

In what follows we first show how rotation solution candidates can be recovered from 7 and 6 matched feature points. Given a rotation solution candidate, it is then shown how the associated translation vector and depths are recovered. Lastly, we show how the unique solution can be distinguished by discarding the physically infeasible solution candidates. The Matlab implementation of QuEst is accessible at <https://goo.gl/QH5qhw>.

We will not discuss why the pose estimation problem has always more than one *mathematically* feasible solution (e.g., 2 for general points and 4 for coplanar points) since these results are well-known. Interested readers are referred to [19] for further discussion and mathematical proofs on the number of solutions.

##### A. Recovering Rotation From 7 Points

Consider the system of equations (11) for 7 matched feature points. Since 7 points generate  $\binom{7}{3} = 35$  equations, in this case  $\mathbf{A}$  is a  $35 \times 35$  matrix. Due to the mathematical multiplicity of solutions however,  $\mathbf{A}$  cannot be full rank (otherwise, only one solution exists, which is a contradiction).

Let us arrange the entries of  $\mathbf{x} \in \mathbb{R}^{35}$  in (11) such that  $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ , where  $\mathbf{x}_1 \in \mathbb{R}^4$  and  $\mathbf{x}_2 \in \mathbb{R}^{31}$  are defined as

$$\mathbf{x}_1 := \begin{bmatrix} w^4 \\ w^3x \\ w^3y \\ w^3z \end{bmatrix}, \quad \mathbf{x}_2 := \begin{bmatrix} x^3w \\ x^4 \\ x^3y \\ x^3z \\ \vdots \end{bmatrix}. \quad (12)$$

Let  $\mathbf{A} = [\mathbf{A}_1 \ \mathbf{A}_2]$ , where  $\mathbf{A}_1 \in \mathbb{R}^{35 \times 4}$  is the first 4 columns of  $\mathbf{A}$ , and  $\mathbf{A}_2 \in \mathbb{R}^{35 \times 31}$  is the remaining part. Equation (11) is equivalent to

$$\mathbf{A}_1 \mathbf{x}_1 + \mathbf{A}_2 \mathbf{x}_2 = \mathbf{0}, \quad (13)$$

from which, by multiplying the pseudo inverse of  $\mathbf{A}_2$  from the left, we obtain<sup>2</sup>

$$\mathbf{x}_2 = -\mathbf{A}_2^\dagger \mathbf{A}_1 \mathbf{x}_1. \quad (14)$$

Let  $\bar{\mathbf{B}} := -\mathbf{A}_2^\dagger \mathbf{A}_1 \in \mathbb{R}^{31 \times 4}$ . The first 4 rows of (14) imply

$$\begin{bmatrix} x^3w \\ x^4 \\ x^3y \\ x^3z \end{bmatrix} = \mathbf{B} \begin{bmatrix} w^4 \\ w^3x \\ w^3y \\ w^3z \end{bmatrix}, \quad (15)$$

<sup>2</sup>For a general point configuration  $\mathbf{A}_2$  has rank 31, and thus  $\mathbf{A}_2^\dagger \mathbf{A}_2 = \mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix.

where  $\mathbf{B} \in \mathbb{R}^{4 \times 4}$  is the matrix consisting of the first 4 rows of  $\bar{\mathbf{B}}$ . By factoring  $x^3$  from the left hand side vector and  $w^3$  from the right hand side vector of (15) we get

$$\frac{x^3}{w^3} \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix} = \mathbf{B} \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix}, \quad (16)$$

which is an eigenvalue problem of the form  $\lambda \mathbf{v} = \mathbf{B}\mathbf{v}$ , with  $\lambda = \frac{x^3}{w^3}$  and  $\mathbf{v} = [w \ x \ y \ z]^\top$ . Hence, 4 solution candidates are found by calculating the (unit norm) eigenvectors of  $\mathbf{B}$ .

We should mention that the choice of  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are somewhat arbitrary. For example, we could have chosen  $\mathbf{x}_2$  as  $\mathbf{x}_2 := [y^3w \ y^3x \ y^4 \ y^3z \ \dots]^\top$ , and derive a similar eigenvalue problem with  $\lambda = \frac{y^3}{w^3}$ . We will later use this fact to distinguish the unique solution.

##### B. Recovering Rotation From 6 Points

Consider equation (11) for 6 feature points. Since 6 feature points generate  $\binom{6}{3} = 20$  equations, in this case  $\mathbf{A}$  is a  $20 \times 35$  full rank matrix.

We split  $\mathbf{x} \in \mathbb{R}^{35}$  into two vectors  $\mathbf{x}_1 \in \mathbb{R}^{20}$  and  $\mathbf{x}_2 \in \mathbb{R}^{15}$ , where  $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ ,  $\mathbf{x}_1$  is the vector of all monomials that contain a power of  $w$  (e.g.,  $w^4, w^3x, w^3y, \dots, wy^3, wz^3$ ), and  $\mathbf{x}_2$  consists of the rest of the monomials (e.g.,  $x^4, x^3y, x^2y^2, \dots, yz^3, z^4$ ). Let  $\mathbf{A} = [\mathbf{A}_1 \ \mathbf{A}_2]$ , where  $\mathbf{A}_1 \in \mathbb{R}^{20 \times 20}$  consists of the first 20 columns of  $\mathbf{A}$ , and  $\mathbf{A}_2 \in \mathbb{R}^{20 \times 15}$  is the remaining part. Equation (11) is equivalent to

$$\mathbf{A}_1 \mathbf{x}_1 + \mathbf{A}_2 \mathbf{x}_2 = \mathbf{0}, \quad (17)$$

from which, by multiplying the pseudo inverse of  $\mathbf{A}_2$ , we obtain

$$\mathbf{x}_2 = -\mathbf{A}_2^\dagger \mathbf{A}_1 \mathbf{x}_1. \quad (18)$$

Since  $\mathbf{x}_1$  consists of monomials that have at least one power of  $w$ , we can factor out  $w$  and represent the remaining vector by  $\mathbf{v} \in \mathbb{R}^{20}$ , i.e.,  $\mathbf{v} = \frac{1}{w} \mathbf{x}_1$ . Thus, (18) can be written as

$$\mathbf{x}_2 = w \bar{\mathbf{B}} \mathbf{v}, \quad (19)$$

where  $\bar{\mathbf{B}} := -\mathbf{A}_2^\dagger \mathbf{A}_1 \in \mathbb{R}^{15 \times 20}$ .

Equation (19) allows us to construct an eigenvalue problem of the form  $\lambda \mathbf{v} = \mathbf{B}\mathbf{v}$ , with  $\mathbf{B} \in \mathbb{R}^{20 \times 20}$ . Indeed, let us choose  $\lambda = \frac{x}{w}$ , and consider the eigenvalue problem

$$\mathbf{x}\mathbf{v} = w \mathbf{B}\mathbf{v}. \quad (20)$$

The entries of vector  $\mathbf{x}\mathbf{v}$  either belong to  $\mathbf{x}_2$  or  $\mathbf{x}_1$ . For entries that belong to  $\mathbf{x}_2$ , the associated rows of  $\mathbf{B}$  are chosen from the corresponding rows of  $\bar{\mathbf{B}}$  in (19). For entries that belong to  $\mathbf{x}_1$ , rows of  $\mathbf{B}$  are chosen as  $[0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$ . The following example illustrates this procedure.

**Example 4.** Suppose entries of  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are arranged as

$$\mathbf{x}_1 = \begin{bmatrix} w^3x \\ w^2x^2 \\ wx^3 \\ \vdots \\ wz^3 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} z^4 \\ xz^3 \\ yz^3 \\ \vdots \\ x^4 \end{bmatrix}, \quad \mathbf{v} = \frac{1}{w} \mathbf{x}_1 = \begin{bmatrix} w^2x \\ wx^2 \\ x^3 \\ \vdots \\ z^3 \end{bmatrix}, \quad (21)$$

and assume that from the feature point coordinates we have derived (19) as

$$\mathbf{x}_2 = \begin{bmatrix} z^4 \\ xz^3 \\ yz^3 \\ \vdots \\ x^4 \end{bmatrix} = w \begin{bmatrix} -0.1 & 5.2 & 22.9 & \dots & 14 \\ 0.1 & -6.7 & -4.4 & \dots & -15 \\ 0.0 & -4.7 & -1.1 & \dots & -11 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -0.1 & -0.4 & 26.2 & \dots & 46 \end{bmatrix} \mathbf{v}. \quad (22)$$

From (22) we can construct the eigenvalue problem (20) as

$$x\mathbf{v} = \begin{bmatrix} w^2x^2 \\ wx^3 \\ x^4 \\ \vdots \\ xz^3 \end{bmatrix} = w \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ -0.1 & -0.4 & 26.2 & \dots & 46 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0.1 & -6.7 & -4.4 & \dots & -15 \end{bmatrix} \mathbf{v}, \quad (23)$$

where the first two entries of  $x\mathbf{v}$  belong to  $\mathbf{x}_1 = w\mathbf{v}$ , and hence their associated rows in  $\mathbf{B}$  consist of zeros except for a single one entry. The third and last entries of  $x\mathbf{v}$  belong to  $\mathbf{x}_2$ , and their associated rows come from  $\bar{\mathbf{B}}$  in (22).

Once the eigenvalue problem (20) is constructed, 20 solution candidates for  $\mathbf{v}$  are derived by computing the eigenvectors of  $\mathbf{B}$ . For each solution candidate,  $w, x, y, z$  are found by calculating the third root of the  $w^3, x^3, y^3, z^3$  entries in  $\mathbf{v}$ . The recovered solution can be normalized to meet the unit norm constraint  $w^2 + x^2 + y^2 + z^2 = 1$ . Notice that by choosing  $\mathbf{x}_1$  or  $\lambda$  differently (e.g.,  $\lambda = \frac{y}{w}$ ) it is possible to derive different eigenvalue problems of the form (19).

### C. Recovering Translation and Depths

Once quaternion elements  $w, x, y, z$  are recovered, the corresponding rotation matrix  $\mathbf{R}$  is given by (5). Having  $\mathbf{R}$ , the rigid motion constraint  $u\mathbf{R}\mathbf{m} + \mathbf{t} = v\mathbf{n}$  can now be written for all matched feature points, and stacked into the matrix-vector form

$$\underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{R}\mathbf{m}_1 & -\mathbf{n}_1 & 0 & 0 & \dots & 0 & 0 \\ \mathbf{I} & 0 & 0 & \mathbf{R}\mathbf{m}_2 & -\mathbf{n}_2 & & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{I} & 0 & 0 & 0 & 0 & \dots & \mathbf{R}\mathbf{m}_k & -\mathbf{n}_k \end{bmatrix}}_{\mathbf{C}} \underbrace{\begin{bmatrix} \mathbf{t} \\ u_1 \\ v_1 \\ u_2 \\ v_2 \\ \vdots \\ u_k \\ v_k \end{bmatrix}}_{\mathbf{y}} = \mathbf{0} \quad (24)$$

where  $\mathbf{I} \in \mathbb{R}^{3 \times 3}$  is the identity matrix,  $k$  is the number feature points,  $\mathbf{C} \in \mathbb{R}^{3k \times 2k+3}$ , and  $\mathbf{y} \in \mathbb{R}^{2k+3}$ . Equation (24) implies that  $\mathbf{y}$  is in the null space of  $\mathbf{C}$ . Thus,  $\mathbf{y}$  can be found by calculating the rightmost singular vector of  $\mathbf{C}$  (i.e., eigenvector of  $\mathbf{C}^T\mathbf{C}$  corresponding to the zero eigenvalue). Notice that  $\mathbf{y}$  consists of the translation vector and feature point depths. Therefore, these parameters are recovered simultaneously and with a common scale factor.

### D. The Unique Solution

Before we proceed with finding the unique solution, we need to briefly talk about the *critical surfaces*. Critical surfaces are special configurations of 3D points in the space for which one cannot distinguish a unique solution. (In this case the problem always has more than one physically

realizable solution.) Perhaps the most important and practical example of such surfaces is when all 3D points lie on a plane, i.e., *coplanar* points. Coplanar points are abundant in aerial images (due to large distance of points from the camera) or images of the man-made environments (due to points lying on walls, floor, etc.). In what follows we will discuss the case of general and coplanar points separately.

Since feature points that are reflected with respect to the origin of the camera frame produce the same image, to detect the physically infeasible solutions one should check the *chirality*. Solution with the wrong chirality correspond to points that are behind the camera, and therefore have negative depths. Hence, physically infeasible solution candidates can be detected and discarded after recovering the depths.

1) *General points*: As discussed at the end of Sections IV-A and IV-B, by choosing other values for  $\lambda$  (e.g.  $\lambda = \frac{y}{w}, \frac{z}{w}$ ) similar eigenvalue problems of the form (16) and (20) can be derived. Since the correct solution must satisfy the eigenvalue problem regardless of the chosen  $\lambda$ , solution candidates that do not satisfy  $\lambda v = \mathbf{B}v$  for all values of  $\lambda$  can be discarded. In this case, two mathematically feasible solutions remain, from which the unique solution is determined by checking the chirality.

2) *Coplanar points*: When points are coplanar (or more generally lie on critical surfaces), four mathematically feasible solutions remain after checking  $\lambda v = \mathbf{B}v$  for different values of  $\lambda$ . Two of these solutions can be discarded by checking the chirality. To determine the solution uniquely further information is required (e.g., a third view or the normal vector to the plane).

Although in theory the methods discussed above can eliminate infeasible solution candidates, in practice pixelization noise and matching imperfections may result in choosing the wrong solution. For instance, in applications where *parallax* is small (i.e., translation between the views is much smaller than the average distance of the 3D points to the camera), recovering the depths becomes an ill-conditioned problem. Thus, checking the chirality may result in the wrong conclusion. Furthermore, with noise, coplanar points may appear as if they are at general positions and instead of two only one solution is returned. It is therefore recommended to always keep the best four solution candidates, and use the Random Sample Consensus (RANSAC) algorithm to find the unique solution.

## V. NOISE AND TIME BENCHMARKS

We benchmark our proposed algorithm against some of the most well-known algorithms such as the essential matrix 8-point [10], 7-point, and 6-point algorithms [19], and the Euclidean homography algorithm [20]. For the first three algorithms the Stewenius's implementation [21], and for the latter Hartley's implementation in Matlab are used. We refer to the algorithms presented in this paper for 6 and 7 points respectively as QuEst 6 and QuEst 7.

Each Monte Carlo simulation consists of eight randomly generated 3D points with uniform distribution inside a rectangular parallelepiped in front of the camera at the initial

location. The camera is moved to a second location by a random translation and rotation quaternion, with uniform distribution within a bounded box of  $\mathbb{R}^3$  and on the 3-sphere, respectively. Coordinates of the feature points on the image plane are computed by projecting the 3D points on the image planes. Each algorithm is provided with the minimum number of points it requires to compute the pose.

### A. Noise Benchmarks

To evaluate the performance under noise, Gaussian noise with zero mean and standard deviation ranging from 0 to 10 pixels is added to all image coordinates. The noise standard deviation is increased by 0.1 pixel increments, and for each noise increment 100 simulations are generated. As mentioned in Section IV-D, the chirality condition is sensitive to noise, so to avoid choosing the wrong solution, the solution candidate that is closest to the ground truth is chosen as the best pose estimate for each algorithm. The estimation error for rotation is defined by

$$\rho(\mathbf{q}, \mathbf{q}^*) = \frac{1}{\pi} \arccos(\text{dot}(\mathbf{q}, \mathbf{q}^*)) \in [0, 1] \quad (25)$$

where  $\mathbf{q} = [w \ x \ y \ z]^T$  is the rotation estimated from the noisy images, and  $\mathbf{q}^*$  is the ground truth rotation in quaternions. Note that (25) defines a metric on the rotation quaternion space [22]. Similarly, the estimation error for translation is defined by

$$\rho(\mathbf{t}_n, \mathbf{t}_n^*) = \frac{1}{\pi} \arccos(\text{dot}(\mathbf{t}_n, \mathbf{t}_n^*)) \in [0, 1] \quad (26)$$

where  $\mathbf{t}_n$  and  $\mathbf{t}_n^*$  are the estimated translation vector and the ground truth, respectively, normalized to have unit norm (because the magnitude of the recovered translation vector can vary depending on the algorithm).

Figure 3 shows the mean rotation and translation estimation errors at different noise standard deviations for all algorithms. Since the feature points are randomly generated and are generally non-coplanar, the homography algorithm, which only works for coplanar points, fails to correctly estimate the pose. Essential matrix based algorithms however are not affected. QuEst 6 and 7-point algorithms have the best performance for rotation, while QuEst 6 and QuEst 7 have the best performance for translation estimates. Unlike the 7-point algorithm, translation estimated by QuEst 7 benefits from the extra points and is comparable to QuEst 6. We should mention that for large noise standard deviations (e.g., 10 pixels) the results can be interpreted as how robust an algorithm is to incorrectly matched feature points.

To analyze the performance when points are on critical surfaces, the previous analysis is repeated for coplanar points, where the points are chosen randomly on a bounded plane with uniform distribution. The mean of the rotation and translation estimation errors for noise standard deviation varying from 0 to 2 is shown in Fig. 4. As can be seen from the figure, homography shows the best noise resilience when the standard deviation is small (approximately 1 to 1.5 pixels). This is because the homography algorithm is specifically designed to recover the pose when points are

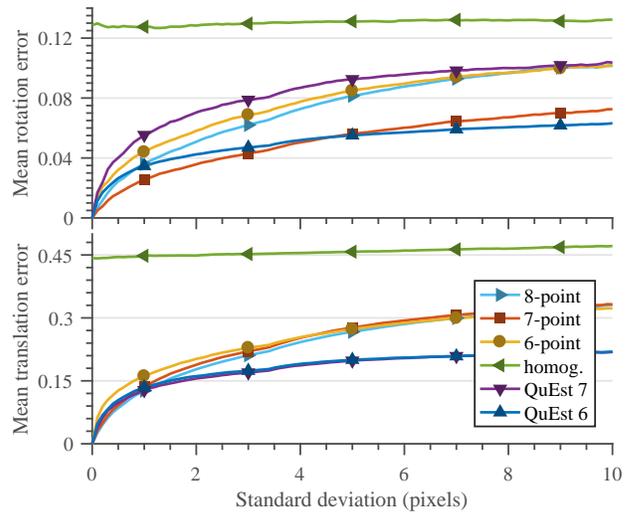


Fig. 3. Comparison under Gaussian noise on feature point coordinates when points are in general 3D configuration.

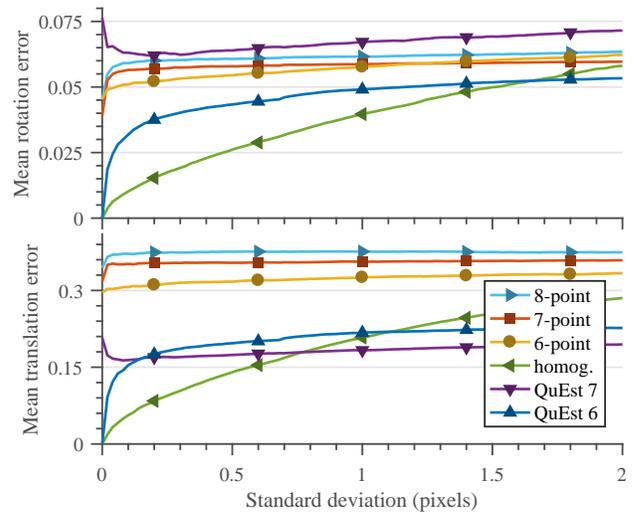


Fig. 4. Comparison under Gaussian noise on feature point coordinates when points are in coplanar configuration.

coplanar. QuEst 6 has the next best estimation accuracy. When points are on critical surfaces matrix  $\mathbf{A}_2$  in (13) loses rank and becomes rank 27. Hence, multiplication by  $\mathbf{A}_2^\dagger$  will not result in (13), and QuEst 7 fails to recover the pose. On the other hand,  $\mathbf{A}_2$  used for QuEst 6 in (17) remains full rank due to having smaller dimensions. Lastly, none of the algorithms that are based on the essential matrix can recover the pose in this case, regardless of the number of points used in the algorithm or the magnitude of noise (see [19] for further explanation).

In conclusion, QuEst 6 shows the best performance since the pose is estimated correctly regardless of the 3D point configuration, and the estimation is robust to noise and outliers.

### B. Time Benchmarks

Table II lists the average execution time of all algorithms in milliseconds, where 1000 Monte Carlo simulations with

TABLE II. Average execution time of essential matrix based algorithms, homography algorithm, and QuEst algorithms.

Algorithm	Average execution time (ms)
Essential 8 point	0.0440
Essential 7 point	0.0288
Essential 6 point	0.0829
Euclidean homography	0.0475
QuEst 7 point	0.5564
QuEst 6 point	0.7728

both coplanar and general points and various noise magnitudes are used to generate the results. All algorithms are implemented as Mex files in Matlab, and tested on the same platform with Intel’s 4th Gen i-7 CPU. Algorithms that use homography or essential matrix have smaller execution time since fewer operations are needed to estimate these matrices. QuEst has a larger execution time since the rotation and translation are recovered independently. We should emphasize that although QuEst is not as fast as other algorithms, it is fast enough to be used with RANSAC in real-time applications. Furthermore, since QuEst recovers the pose regardless of the 3D point configuration, no prior effort is required to detect the coplanarity and choose the appropriate algorithm correspondingly.

## VI. REAL WORLD PERFORMANCE

To further assess the accuracy and robustness of QuEst, images from four real world datasets are used to estimate the pose and compare the results with the ground truth that is provided by the dataset. Datasets used for the comparison are ICL [23], KITTI [24], NAIST [25], and TUM [26], where SURF image feature points are extracted and matched between two consecutive keyframes. Each algorithm is provided with the minimum number of points it requires to estimate the pose, where points with the highest matching score are chosen. To test the robustness of the algorithms, RANSAC is not used, and the provided set of matched points can occasionally have outliers.

The ICL-NUIM dataset consists of computer-generated renderings of 3D scenes. These artificial images are accompanied by exact ground truth information, and their main purpose is to benchmark 3D reconstruction algorithms. The KITTI and TUM datasets were created to evaluate SLAM algorithms. The KITTI dataset consists of outdoor stereo image sequences taken from a moving vehicle. The images are accompanied by LIDAR, IMU and GPS measurements, so full pose information is provided in the left camera frame. The TUM dataset is comprised of indoor images as well as point depth information from a RGB-D camera. The RGB-D camera poses were recorded by using an optical tracking system, with millimeter-level accuracy. The NAIST campus sequences are part of an Augmented Reality benchmark called TrakMark. The ground truth files for these sequences were created by solving a PnP problem from known 3D world point coordinates, acquired using high

precision surveying equipment. The images come from a handheld camera, with the operator walking while capturing some sequences and running for others.

Figure 5 shows the rotation and translation estimation errors associated to each algorithm for all datasets. As shown in the figure, most algorithms perform well on the ICL benchmark, due to its lack of image noise. The translation error for all algorithms is the lowest for the KITTI sequences. This can be explained by the camera motion, since most of the time the camera is moving forward, in the direction of its z-axis. The sequences labeled as NAIST were challenging for every algorithm, since the camera motion is erratic at times and changes quickly. The TUM sequences, while generated by a handheld sensor, do not contain motions as abrupt as those in NAIST. The performance of all algorithms on the TUM sequences is better than for the NAIST sequences, but below their performance on the ICL image sets.

Since the chosen feature points may not be coplanar, the homography algorithm does not return the correct solution in general. The pose estimated by the 6-point and 7-point algorithms has smaller median error compared to the 8-point algorithm. The translation estimated by QuEst 7 has the smallest median error due to the additional point. The rotation estimated by QuEst 6 shows the best performance due to having smaller median errors, 0.25 and 0.75 quartiles, and outliers for all datasets.

## VII. CONCLUSION AND FUTURE WORK

By using quaternion representation of rotation, we formulated the camera pose estimation problem and presented the QuEst algorithm to recover the relative pose between two camera views. Unlike the existing homography or essential matrix based methods, QuEst decouples the rotation and translation estimation, and recovers the pose correctly for both cases of general and coplanar points. QuEst can be used to initialize the bundle adjustment algorithm in applications such as SLAM without needing to resort to heuristic methods to detect the coplanarity of the points. Using both simulated and real world images, we demonstrated the estimation accuracy and robustness of QuEst in comparison to the commonly used algorithms. We have made the Matlab implementation of QuEst available online and free. Future work includes the online release of the C++ implementation of QuEst with RANSAC.

## REFERENCES

- [1] S. Agarwal, N. Snavely, I. Simon, S. M. Seitz, and R. Szeliski, “Building rome in a day,” in *IEEE International Conference on Computer Vision*, Sept 2009, pp. 72–79.
- [2] N. Gans, G. Hu, and W. E. Dixon, “Image based state estimation,” in *Encyclopedia of Complexity and Systems Science*. Springer, 2009, pp. 4751–4776.
- [3] L. Kneip, P. Furgale, and R. Siegwart, “Using multi-camera systems in robotics: Efficient solutions to the NPnP problem,” in *International Conference on Robotics and Automation*. IEEE, 2013, pp. 3770–3776.
- [4] D. Tick, A. C. Satici, J. Shen, and N. Gans, “Tracking control of mobile robots localized via chained fusion of discrete and continuous epipolar geometry, IMU and odometry,” *IEEE transactions on cybernetics*, vol. 43, no. 4, pp. 1237–1250, 2013.

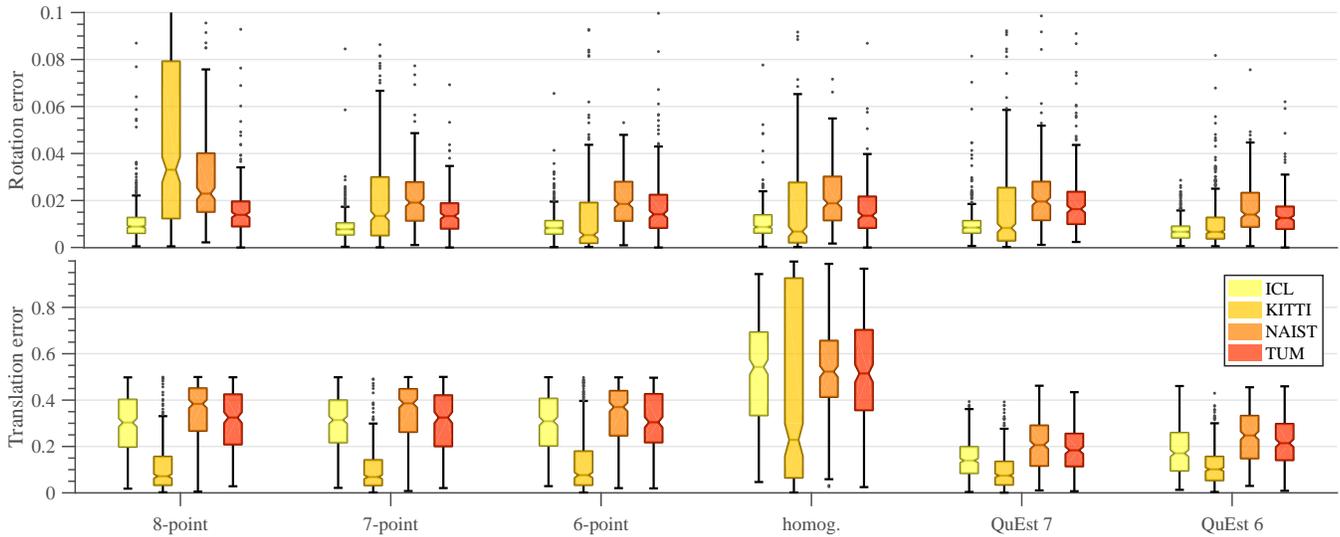


Fig. 5. Rotation and translation estimation error associated to six algorithms tested on image datasets ICL, KITTI, NAIST, and TUM.

- [5] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "ORB-SLAM: a versatile and accurate monocular SLAM system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [6] J. Schneider, C. Eling, L. Klingbeil, H. Kuhlmann, C. Stachniss, et al., "Fast and effective online pose estimation and mapping for UAVs," in *International Conference on Robotics and Automation*. IEEE, 2016, pp. 4784–4791.
- [7] E. Marchand, H. Uchiyama, and F. Spindler, "Pose estimation for augmented reality: a hands-on survey," *IEEE transactions on visualization and computer graphics*, vol. 22, no. 12, pp. 2633–2651, 2016.
- [8] H. Bay, T. Tuytelaars, and L. Van Gool, "SURF: Speeded up robust features," in *European conference on computer vision*. Springer, 2006, pp. 404–417.
- [9] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: an efficient alternative to SIFT or SURF," in *International conference on computer vision*. IEEE, 2011, pp. 2564–2571.
- [10] H. C. Longuet-Higgins, "A computer algorithm for reconstructing a scene from two projections," *Nature*, vol. 293, no. 5828, pp. 133–135, 1981.
- [11] M. Demazure, "Sur deux problemes de reconstruction," INRIA, Tech. Rep. RR-0882, July 1988.
- [12] O. Faugeras, *Three-dimensional computer vision: a geometric viewpoint*. MIT Press, 1993.
- [13] E. Malis and F. Chaumette, "2 1/2 D visual servoing with respect to unknown objects through a new estimation scheme of camera displacement," *International Journal of Computer Vision*, vol. 37, no. 1, pp. 79–97, 2000.
- [14] E. Kruppa, "Zur Ermittlung eines Objektes aus zwei Perspektiven mit innerer Orientierung," *Sitzungsberichte der Mathematisch Naturwissenschaftlichen Kaiserlichen Akademie der Wissenschaften*, vol. 122, pp. 1939–1948, 1913.
- [15] O. D. Faugeras and S. Maybank, "Motion from point matches: multiplicity of solutions," in *Workshop on Visual Motion*, 1989, pp. 248–255.
- [16] J. Philip, "A non-iterative algorithm for determining all essential matrices corresponding to five point pairs," *The Photogrammetric Record*, vol. 15, no. 88, pp. 589–599, 1996.
- [17] K. Fathian and N. R. Gans, "A new approach for solving the five-point relative pose problem for vision-based estimation and control," in *American Control Conference*. IEEE, 2014, pp. 103–109.
- [18] J.-Y. Bouguet, "Camera calibration toolbox," [https://www.vision.caltech.edu/bouguetj/calib\\_doc/](https://www.vision.caltech.edu/bouguetj/calib_doc/), 2015.
- [19] Y. Ma, S. Soatto, J. Kosecka, and S. S. Sastry, *An Invitation to 3-D Vision*. Springer, 2003.
- [20] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge University Press, 2004.
- [21] H. Stewenius, "Eight-point, seven-point, and six-point solvers," <http://www.vis.uky.edu/~stewe/FIVEPOINT>, 2006.
- [22] D. Q. Huynh, "Metrics for 3D rotations: Comparison and analysis," *Journal of Mathematical Imaging and Vision*, vol. 35, no. 2, pp. 155–164, 2009.
- [23] A. Handa, T. Whelan, J. McDonald, and A. Davison, "A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM," in *IEEE Intl. Conf. on Robotics and Automation*, May 2014.
- [24] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the KITTI vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition*, 2012.
- [25] H. Tamura and H. Kato, "Proposal of international voluntary activities on establishing benchmark test schemes for AR/MR geometric registration and tracking methods," in *International Symposium on Mixed and Augmented Reality*. IEEE, 2009, pp. 233–236.
- [26] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of RGB-D SLAM systems," in *International Conference on Intelligent Robot Systems*, Oct. 2012.