

On Computing The Perspective Transformation Matrix and Camera Parameters[†]

T. N. Tan, G. D. Sullivan and K. D. Baker

Department of Computer Science
University of Reading, Berkshire RG6 2AY, UK
Email: T.Tan@uk.ac.reading

Abstract

Camera calibration often entails the computation of the perspective transformation matrix. Conventionally, the matrix has been calculated by the standard linear least squares technique. Recently, Faugeras and Toscani have criticised the conventional approach for producing unsatisfactory, even "absurd", solutions, and have proposed an alternative approach. It is shown in this paper that their criticism of the conventional approach is misplaced and misleading. Experimental results demonstrate that Faugeras and Toscani's approach has no advantage over the conventional approach from the practical point of view. In fact, the latter is shown to be superior both in noise robustness and in computational cost. The paper also reports a method to resolve the possible sign ambiguities in the camera parameters computed by existing algorithms.

1 Introduction

Camera calibration is a classical issue in close-range photogrammetry and a prerequisite for many computer vision applications. It has received considerable attention from the photogrammetry community [1] and more recently from the computer vision community [2]. A good review of the existing techniques may be found in Tsai [2].

A popular camera calibration paradigm achieves camera calibration by first computing the so-called *perspective transformation matrix* (PTM) [3], and then decomposing the matrix into intrinsic and extrinsic camera parameters [4]. This class of algorithms usually assumes a pinhole camera model, with no lens distortion, and provides a closed-form solution. Two recent variants are due to Faugeras and Toscani [5-6], and Puget and Skordas [7].

The PTM relates the 3D world coordinates of points to their 2D image coordinates. It is common practice to seek many such correspondences, so that the PTM is strongly overconstrained. The solution has then conventionally been computed by the standard linear least squares (LLS) technique by setting one of the matrix elements to unity [8]. In their recent papers [5-6], Faugeras and Toscani argue that arbitrarily setting one of the matrix elements to unity causes the

†. This work was carried out as part of the ESPRIT project P2152 (VIEWS).

standard LLS approach to yield "... a solution which is absurd since the intrinsic parameters depend upon the choice of the world coordinate system" [5, p.240]. They then proposed an eigenvector solution based on a physical constraint derived from the elements of the PTM. Their criticism of the conventional approach has apparently been unquestioned and widely accepted by the computer vision community, as evidenced by the frequent citation of their papers [9-11], and the adoption of their eigenvector approach [7].

In this paper, we present a second look at the two approaches to the estimation of the PTM. We argue that Faugeras and Toscani's criticism of the LLS approach is misplaced since it is based on an inappropriate interpretation of the derived perspective transformation matrix. Both analytical and experimental results show that the use of the conventional approach is well justified, and may be preferable to the eigenvector approach in practical applications.

The second step in a PTM-based camera calibration algorithm is to derive camera parameters from the estimated PTM [4-6]. Closed-form solutions have been presented by Faugeras and Toscani [5-6], and Puget and Skordas [7] among many others. However, some of the parameters calculated by existing algorithms suffer from a possible sign ambiguity, the resolution of which was not clearly stated in these papers. Formulae are given in this paper for resolving the sign ambiguity so as to ensure the orthonormality of the rotation matrix.

2 Estimation of the perspective transformation matrix

We first discuss PTM estimation. After briefly describing the LLS approach, we examine Faugeras and Toscani's criticism and argue that it is ill-founded. We then present experimental results to support our analysis and to compare the performance of the LLS, and Faugeras and Toscani's eigenvector approach.

2.1 The conventional linear least squares approach

Similar notation to those used in [5-6] are adopted here for the sake of easy cross-reference. The camera is assumed to be a pinhole camera with perspective projection and with no lens distortion as in [5-6]. The 3D world coordinates \vec{P} and the 2D image coordinates \vec{p} of a control point are related to each other by

$$\vec{p} = M\vec{P} \quad (1)$$

where

$$M = \begin{bmatrix} l_{11} & l_{12} & l_{13} & l_{14} \\ l_{21} & l_{22} & l_{23} & l_{24} \\ l_{31} & l_{32} & l_{33} & l_{34} \end{bmatrix} = \begin{bmatrix} l_1 & l_{14} \\ l_2 & l_{24} \\ l_3 & l_{34} \end{bmatrix} \quad (2)$$

is the PTM to be estimated. For each control point with known world coordinates $X_i = [x_i \ y_i \ z_i]^T$ and image coordinates $x_i = [u_i \ v_i]^T$, Eqn. (1) yields the

following two linear equations in the 12 unknown elements of M [5-6]:

$$\begin{cases} l_1 X_i - u_i l_3 X_i + l_{14} - u_i l_{34} = 0 \\ l_2 X_i - v_i l_3 X_i + l_{24} - v_i l_{34} = 0 \end{cases} \quad (3)$$

Thus for N control points with known world and image coordinates, one obtains a set of $2N$ linear homogeneous equations:

$$AL = \mathbf{0}_{2N} \quad (4)$$

where A is the known $2N \times 12$ coefficient matrix; $L = [l_1 \ l_{14} \ l_2 \ l_{24} \ l_3 \ l_{34}]^T$ the unknown 12×1 vector; and $\mathbf{0}_{2N}$ a $2N \times 1$ zero vector.

This raises the point under debate: because of the homogeneous nature of (4), the PTM M can only be determined up to a non-zero scale factor, that is, if M is a valid solution of (4), then for any non-zero scale α , αM is also a valid solution. The scalar α is irrelevant to the task although it may have an effect on the numerical accuracy of the PTM.

The LLS approach solves L and M from (4) by letting $l_{34} = 1$, and performing a pseudo-inversion of the coefficient matrix [8]:

$$L = -(C^T C)^{-1} C^T B \quad (5)$$

where C comprises the first 11 columns of A , and B is the last column of A .

2.2 A critique of Faugeras and Toscani's criticism of the conventional approach

To investigate how the PTM computed by the standard LLS approach changes in a new world coordinate system (WCS), Faugeras and Toscani [5-6] suggested transforming the world coordinates of the given N control points by a rotation R followed by a translation t , then computing the PTM in the new WCS, and finally comparing the two PTMs obtained before and after displacing the WCS. Let the old and new world coordinates X_{old} and X_{new} be related to each other by

$$X_{old} = R X_{new} + t \quad (6)$$

Let the PTM corresponding to the old and new WCS be M and M' respectively. Then it can easily be shown based on (1) and (6) that M and M' are related to each other by (assuming M and M' have appropriate scales):

$$M' = \begin{bmatrix} l_1 R & l_1 t + l_{14} \\ l_2 R & l_2 t + l_{24} \\ l_3 R & l_3 t + l_{34} \end{bmatrix} = M \begin{bmatrix} R & t \\ 0_3^T & 1 \end{bmatrix} = MU \quad (7)$$

where $U = \begin{bmatrix} R & t \\ 0_3^T & 1 \end{bmatrix}$. As the 3D coordinates have been changed, (4) becomes

$$\mathbf{A}'\mathbf{L} = \mathbf{0}_{2N} \quad (8)$$

where \mathbf{A}' is the new coefficient matrix depending on the new 3D world coordinates and the unchanged 2D image coordinates.

Let \mathbf{K} and \mathbf{K}' be respectively the PTMs computed from (4) and (8) by the LLS approach. Then based on (7), Faugeras and Toscani [5-6] claim that

“we expect them [\mathbf{K} and \mathbf{K}'] to verify:

$$\mathbf{K}' = \mathbf{K}\mathbf{U} \quad (9)$$

This is not the case if we use the constraint $l_{34} = 1$. In particular, the intrinsic parameters will depend on the choice of the world coordinate system, which is clearly not satisfactory.” [6, p.17].

But the global scale of the PTM computed by the LLS approach is necessarily uncertain, so that in general, one has $\mathbf{M}' = s_1 \mathbf{K}'$ and $\mathbf{M} = s_2 \mathbf{K}$. Therefore \mathbf{K} and \mathbf{K}' are expected to satisfy

$$s_1 \mathbf{K}' = s_2 \mathbf{K}\mathbf{U} \quad (10)$$

or equivalently

$$\mathbf{K}' = s \mathbf{K}\mathbf{U} \quad (11)$$

where s_1, s_2 and $s (=s_2/s_1)$ are non-zero scaling factors. It can be seen from (7) that s is given by

$$s = \frac{1}{\mathbf{m}_3 \mathbf{t} + 1} \quad (12)$$

where \mathbf{m}_3 is a 1x3 row vector composed of the first three elements of the last row of \mathbf{K} . The validity of Equations (10) and (11) can easily be verified. In particular, *the intrinsic parameters obtained in the old and new coordinate systems are identical*. From (10) and (11) it is also evident that (9) holds if the two world coordinate systems are related to each other by a pure rotation (i.e., $\mathbf{t} = \mathbf{0}$), or the translation vector \mathbf{t} is orthogonal to \mathbf{m}_3 .

To overcome the “absurd” solutions produced by the LLS approach, Faugeras and Toscani proposed an alternative based on the constraint $\|\mathbf{l}_3\|^2 = 1$ [5-6]. The constraint, as pointed out by Faugeras and Toscani, is invariant to the displacement of the WCS. The PTM can then be obtained by solving an eigenvector problem [5-6].

2.3 The conventional approach vs. Faugeras and Toscani’s approach

So far it has tacitly been assumed that the input data is noise-free. Under more realistic, noisy conditions, neither (9) nor (10) can be expected to be exact, and the intrinsic camera parameters computed by the LLS approach will appear to depend on the choice of the world coordinate system. This is not unusual. It is well-known

that the accuracy of the solutions to many similar problems of camera calibration, stereo, and structure from motion is dependent (sometimes critically) on the imaging geometry. The recovery of the central transformation from the world coordinate system to the camera coordinate system depends in practice on the specific input data used. Furthermore, the pinhole camera model is only an approximation of the imaging process of physical cameras, and there is often an unresolvable ambiguity between different camera parameters. In our view, under noisy conditions, it makes little sense to emphasize the invariance of the intrinsic parameters to the WCS. What is more desirable from the practical point of view is the robustness of the recovered parameters in the presence of noise.

We have therefore compared the performance of the two algorithms under noisy conditions. Monte Carlo simulations were conducted to examine the robustness of the four intrinsic parameters (the piercing point (U_0, V_0) , and the horizontal and vertical scales α_u and α_v) recovered by the two approaches under noisy conditions before and after displacements of the WCS. Using a representative set of known camera parameters (similar to those of a real calibrated outdoor camera), a specified number (=15) of points were randomly generated in a sphere of a given radius. Noise was simulated by perturbing the image coordinates of each point by a random amount $\Delta\epsilon$ uniformly distributed over $[-\Delta E, +\Delta E]$, where ΔE specifies the noise level. To move the WCS, the world coordinates of all points were transformed by a random motion, and the noisy image data was unchanged.

The PTM was computed twice using the LLS approach and twice using Faugeras and Toscani's approach (once in the old WCS and once in the new WCS). The four PTMs obtained were then subjected to the same PTM decomposition technique [6] to calculate the four intrinsic parameters. At each noise level, 500 trials were conducted, and the mean absolute relative errors of the four intrinsic parameters were computed. The simulation results are summarized in Fig.1. As expected, the PTMs computed in the old and new WCS by Faugeras and Toscani's approach always result in identical intrinsic parameters, whereas those by the LLS approach diverge at high noise levels. However, it is important to note that the horizontal and vertical scales computed by the conventional approach are significantly more accurate than those by Faugeras and Toscani's approach, especially under severe noise conditions. The conventional approach, as a whole, appears to be more noise-robust than Faugeras and Toscani's approach.

The two approaches were also used to calibrate a number of real cameras in both indoor and outdoor scenes. An example is given in Fig.2 where an airport scene is shown. 24 control points were set up manually, and 3D coordinates were obtained by means of on-the-spot measurements. In practice, this is very difficult to do accurately, especially when there is only very limited access to the site (in the case in Fig.2, all measurements had to be carried out within one hour). In traffic scenes, it is generally unreasonable to hope for very accurate calibration

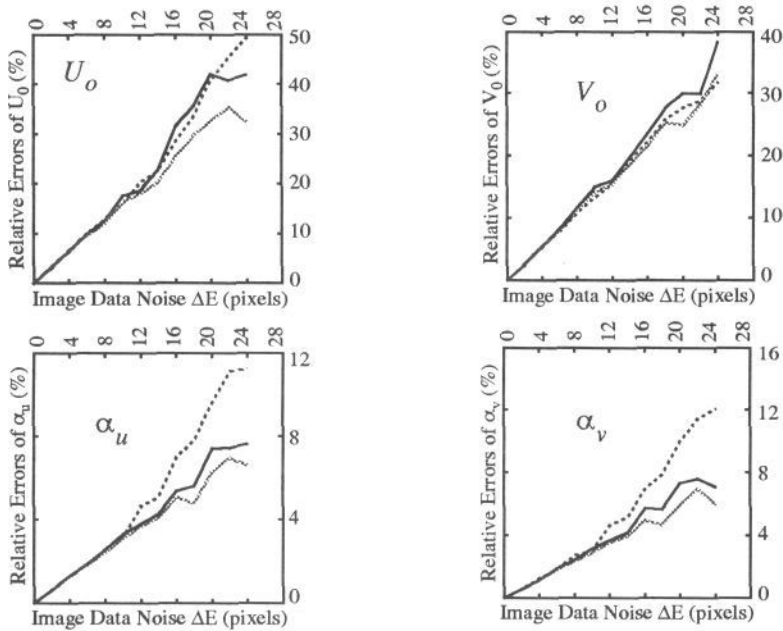


Figure 1: Relative errors of four intrinsic camera parameters. Dotted curve - Faugeras and Toscani's approach both in the old and the new WCS; grey curve - the conventional approach in the old WCS; and dark curve - the conventional approach in the new WCS.

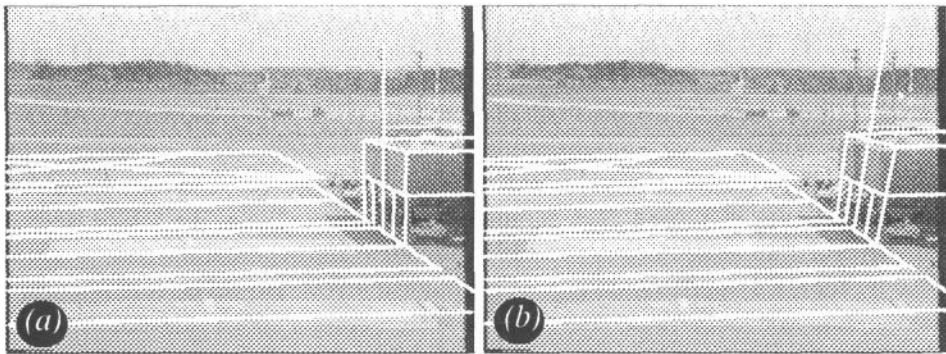


Figure 2: Projection (superimposed on the original image) of a simple scene model when viewed from the camera calibrated by the LLS approach (a) and by Faugeras and Toscani's approach (b).

measurements. Even in a laboratory model of a traffic scene, where we have ad lib access for measurement purposes, errors in the calibration points have proved to be a very significant practical problem. The use of highly accurate calibration grids [2, 5] overcomes this to some extent, but is infeasible outside the laboratory.

Fig.2(a) and (b) show, respectively, the projection (white line segments superimposed on the original airport image) of a simple scene model when viewed

from the camera after calibrated by the two methods. Qualitatively, the conventional approach is seen to perform better than Faugeras and Toscani's approach (compare the projections of the building on the right of the image).

In addition to its performance superiority (at least in the examples shown here), the conventional approach is computationally much simpler than Faugeras and Toscani's approach especially when the number of control points increases. Fig.3 illustrates typical results.

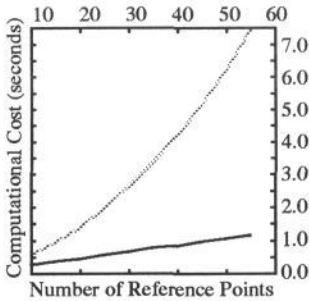


Figure 3: Comparison between the computational cost of the standard linear least squares approach (dark curve), and Faugeras and Toscani's approach (grey curve). All code was written in Pop11 and executed on a Sun 4.

In summary, we have shown in this section that

- The intrinsic camera parameters computed by the conventional approach are reasonably invariant to the choice of the WCS at low noise levels.
- As a whole, the conventional approach seems more noise-robust than Faugeras and Toscani's approach.
- The conventional approach often appears superior to Faugeras and Toscani's approach with real data.
- The conventional approach is computationally much simpler than Faugeras and Toscani's approach especially when a large number of control points are used.

We therefore can conclude that Faugeras and Toscani's criticism of the conventional LLS approach is not generally justified. In particular, Faugeras and Toscani's approach has been shown to offer no advantage over the conventional approach, which, in fact, has important practical benefits.

3 Resolving sign ambiguities in camera parameters

The exact form of M as a function of camera parameters depends on the specific camera model. If the camera is a pinhole camera with no lens distortion, and has the geometry illustrated in Fig.4, then M is given by [6]

$$M = k \begin{bmatrix} \alpha_u r_1 + U_0 r_3 & \alpha_u t_x + U_0 t_z \\ \alpha_v r_2 + V_0 r_3 & \alpha_v t_y + V_0 t_z \\ r_3 & t_z \end{bmatrix} = \begin{bmatrix} l_1 & l_{14} \\ l_2 & l_{24} \\ l_3 & l_{34} \end{bmatrix} \quad (13)$$

where k is a non-zero scale factor, (U_0, V_0) the piercing point, α_u and α_v the

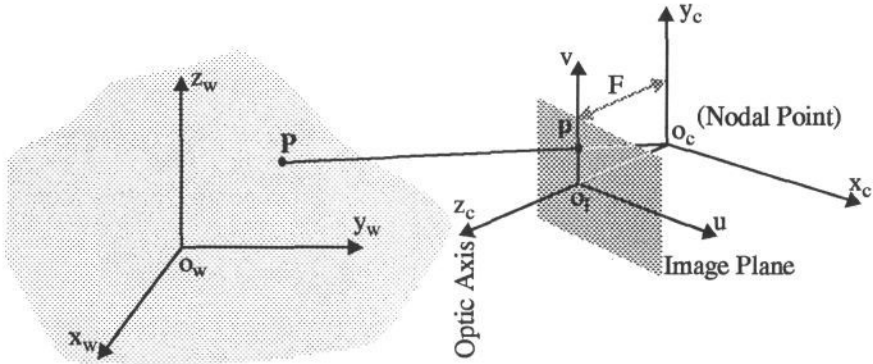


Figure 4: Pinhole camera with perspective projection and no lens distortion.

horizontal and vertical scales, $(t_x \ t_y \ t_z)^T$ the translation vector, and $\mathbf{r}_1, \mathbf{r}_2$ and \mathbf{r}_3 the three row vectors of the 3×3 rotation matrix. U_0, V_0, α_u and α_v are the four intrinsic parameters, and $t_x, t_y, t_z, \mathbf{r}_1, \mathbf{r}_2$ and \mathbf{r}_3 define the six extrinsic parameters.

Closed-form solutions are presented in Puget and Skordas [7] for recovering the 10 camera parameters from the PTM given in (13). They assume that the PTM has been normalized so that $k = 1$, and calculate the four intrinsic parameters based on the following equations:

$$U_0 = \mathbf{l}_1 \cdot \mathbf{l}_3; \quad V_0 = \mathbf{l}_2 \cdot \mathbf{l}_3; \quad \alpha_u = \|\mathbf{l}_1 \times \mathbf{l}_3\|; \quad \alpha_v = \|\mathbf{l}_2 \times \mathbf{l}_3\| \quad (14)$$

where \mathbf{l}_i replaces \mathbf{m}_i used in [7]. On using (14), Puget and Skordas implicitly assume that both scaling factors α_u and α_v are positive. Nevertheless, when the origin of the WCS is in front of the camera and t_z is assumed to be positive (as is usually the case), then only one of the two scaling factors can arbitrarily be chosen to be positive [4]. The other scaling factor can be either positive or negative. This implies that one of the scaling factors calculated by (14) is subject to a sign ambiguity. Assume, for the moment, that α_u is positive. As α_v and \mathbf{r}_2 , and α_v and t_y always appear in product forms, \mathbf{r}_2 and t_y suffer from the same sign ambiguity, and the rotation matrix so obtained may not be orthonormal.

To resolve the sign ambiguity associated with α_v , we first compute U_0, V_0 and α_u as in (14). Then from (13) and using the orthonormality of the rotation matrix, we have

$$\mathbf{r}_3 = \mathbf{l}_3; \quad \mathbf{r}_1 = \frac{1}{\alpha_u} (\mathbf{l}_1 - U_0 \mathbf{r}_3) \quad (15)$$

$$\mathbf{l}_2 \times \mathbf{l}_3 = (\alpha_v \mathbf{r}_2 + V_0 \mathbf{r}_3) \times \mathbf{r}_3 = \alpha_v \mathbf{r}_2 \times \mathbf{r}_3 = \alpha_v \mathbf{r}_1$$

By combining the three equations in (15), we obtain

$$\alpha_v = \frac{1}{\alpha_u} (\mathbf{l}_1 - U_0 \mathbf{l}_3) \cdot (\mathbf{l}_2 \times \mathbf{l}_3) \quad (16)$$

Clearly, α_v computed above suffers from no sign ambiguity. If α_v is chosen to be positive, then a unique solution for α_u can similarly be obtained.

In the camera model adopted so far, the u and v axes of the image plane are assumed to be perfectly perpendicular. To consider a possible non-perpendicularity of the two axes, Faugeras and Toscani [5] suggest modelling the mapping from retina coordinates (u, v) to raster image coordinates (x, y) by:

$$x = a + bu + cv; \quad y = d + ev \quad (17)$$

where a, b, c, d and e represent five intrinsic camera parameters with $1/b$ and $1/e$ being the two scale factors. In this case, the PTM is given by [5]

$$M = k \begin{bmatrix} \frac{1}{b}r_1 - \frac{c}{be}r_2 + gr_3 & \frac{1}{b}t_x - \frac{c}{be}t_y + gt_z \\ \frac{1}{e}r_2 - \frac{d}{e}r_3 & \frac{1}{e}t_y - \frac{d}{e}t_z \\ r_3 & t_z \end{bmatrix} = \begin{bmatrix} l_1 & l_{14} \\ l_2 & l_{24} \\ l_3 & l_{34} \end{bmatrix} \quad (18)$$

where $g = (cd - ae)/be$. Formulae for computing the camera parameters from (18) are given by Faugeras and Toscani [5]. They assume $t_z, e > 0$, and compute b by

$$b = \frac{1}{\sqrt{\frac{l_1 l_1^T}{k^2} - \left(\frac{c}{b}\right)^2 \frac{1}{e^2} - g^2}} \quad (19)$$

where $k, (c/b), e$ and g have been pre-computed. Although they point out that b can be either positive or negative, and that (19) only determines the magnitude of b , it is not discussed in [5] how to resolve the sense of b . It can be seen from (18) that if b is subject to a sign ambiguity, r_1 and t_x are subject to the same ambiguity, and the resultant rotation matrix may be non-orthonormal. To determine a unique b , we calculate k and r_2 as in [5]. Then from the orthonormal properties of the rotation matrix and (18), we have

$$r_2 \cdot (l_3 \times l_1) = r_2 \cdot \left(\frac{c}{be}r_1 + \frac{1}{b}r_2 \right) k^2 = \frac{k^2}{b} \quad (20)$$

i.e.,

$$b = k^2 / (r_2, l_3, l_1) \quad (21)$$

where (\cdot, \cdot, \cdot) denotes the triple scalar product. Clearly, b computed by (21) is subject to no sign ambiguity. If we choose b to be positive, a unique solution for the other scale factor e can similarly be obtained.

4 Conclusions

The perspective transformation matrix has conventionally been computed by the linear least squares technique and the inherent scale ambiguity in homogeneous algebra is resolved by setting one of the unknowns to unity. The conventional approach has recently been criticised by Faugeras and Toscani for producing “absurd” solutions which depend on the (arbitrary) choice of the world coordinate frame. We have shown in this paper that the criticism reflects an over-strict interpretation of the perspective transformation matrix. Experimental results have been presented which justify our objection to Faugeras and Toscani’s criticism of the conventional approach. We have found the approach proposed by Faugeras and Toscani to have no advantage over the conventional approach in practice. In fact the opposite is true: in our experiments, the conventional approach has been found superior both in noise robustness and in computational cost.

It has also been shown that the possible sign ambiguities in the camera parameters computed by existing algorithms can easily be resolved by examining the orthonormality of the rotation matrix. Formulae for obtaining unique camera parameters have been given as minor modifications of the existing algorithms.

5 References

1. H. M. Karara, Close-Range Photogrammetry: Where Are We and Where Are We Heading, *Photogrammetric Eng. Remote Sensing*, vol.51, 1985, pp.537-544.
2. R. Y. Tsai, Synopsis of Recent Progress on Camera Calibration for 3D Machine Vision, in *The Robotics Review*, O. Khatib, et. al., (Eds.), MIT Press, 1989.
3. R. M. Haralick, Using Perspective Transformations in Scene Analysis, *Computer Graphics and Image Processing*, vol.13, 1980, pp.191-221.
4. S. Ganapathy, Decomposition of Transformation Matrices for Robot Vision, *Pattern Recognition Letter*, vol.2, 1984, pp.401-412.
5. O. D. Faugeras and G. Toscani, Camera Calibration for 3D Computer Vision, *Proc. of IEEE Int. Workshop on Machine Vision & Machine Intelligence*, Tokyo, Japan, Feb. 1987, pp.240-247.
6. O. D. Faugeras and G. Toscani, The Calibration Problem for Stereo, *Proc. of IEEE Int. Conf. Computer Vision and Pattern Recognition*, Miami Beach, FL, USA, 1986, pp.15-20.
7. P. Puget and T. Skordas, Calibrating a mobile camera, *Image and Vision Computing*, vol. 8, no. 4, November 1990, pp.341-348.
8. D. H. Ballard and C. M. Brown, *Computer Vision*, New Jersey: Prentice-Hall, 1982.
9. G. Q. Wei and S. D. Ma, Two Plane Camera Calibration: A Unified Model, *Proc. of CVPR'91*, Maui, Hawaii, USA, June 3-6, 1991, pp.133-138.
10. P. Beardsley, D. Murray and A. Zisserman, Camera Calibration Using Multiple Images, *Proc. of ECCV92*, LNCS-588, Springer-Verlag, 1992, pp.312-320.
11. C.-C. Wang, Extrinsic Calibration of a Vision Sensor Mounted on a Robot, *IEEE Trans. Robotics and Automation*, vol.8, 1992, pp.161-175.