

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
ARTIFICIAL INTELLIGENCE LABORATORY
and
CENTER FOR BIOLOGICAL COMPUTATIONAL LEARNING
WHITAKER COLLEGE

A.I. Memo No. 1452
C.B.C.L. Paper No. 90

January, 1994

Algebraic Functions For Recognition

Amnon Shashua

Abstract

In the general case, a trilinear relationship between three perspective views is shown to exist. The *trilinearity* result is shown to be of much practical use in visual recognition by alignment — yielding a direct method that cuts through the computations of camera transformation, scene structure and epipolar geometry. The proof of the central result may be of further interest as it demonstrates certain regularities across homographies of the plane and introduces new view invariants. Experiments on simulated and real image data were conducted, including a comparative analysis with epipolar intersection and the linear combination methods, with results indicating a greater degree of robustness in practice and a higher level of performance in re-projection tasks.

Copyright © Massachusetts Institute of Technology, 1994

This report describes research done within the Center for Biological and Computational Learning in the Department of Brain and Cognitive Sciences, and at the Artificial Intelligence Laboratory. Support for the A.I. Laboratory's artificial intelligence research is provided in part by the Advanced Research Projects Agency of the Department of Defense under Office of Naval Research contract N00014-91-J-4038. Support for the Center's research is provided in part by ONR contracts N00014-91-J-1270 and N00014-92-J-1879; by a grant from the National Science Foundation under contract ASC-9217041 (funds provided by this award include funds from ARPA provided under HPCC); and by a grant from the National Institutes of Health under contract NIH 2-S07-RR07047-26. Additional support is provided by the North Atlantic Treaty Organization, ATR Audio and Visual Perception Research Laboratories, Mitsubishi Electric Corporation, Siemens AG., and Sumitomo Metal Industries. A. Shashua is supported by a McDonnell-Pew postdoctoral fellowship from the department of Brain and Cognitive Sciences.

1 Introduction

We establish a general result about algebraic connections across three perspective views of a 3D scene and demonstrate its application to visual recognition via alignment. We show that, in general, any three perspective views of a scene satisfy a pair of trilinear functions of image coordinates. In the limiting case, when all three views are orthographic, these functions become linear and reduce to the form discovered by [34]. Using the trilinear result one can manipulate views of an object (such as generate novel views from two model views) without recovering scene structure (metric or non-metric), camera transformation, or even the epipolar geometry.

The central results in this paper are contained in Theorems 1 and 2. The first theorem states that the variety of views ψ of a fixed 3D object obtained by an uncalibrated pin-hole camera satisfy a relation of the sort $F(\psi, \psi_1, \psi_2) = 0$, where ψ_1, ψ_2 are two arbitrary views of the object, and F has a special trilinear form. The coefficients of F can be recovered linearly without establishing first the epipolar geometry, 3D structure of the object, or camera motion. The auxiliary Lemmas required for the proof of Theorem 1 may be of interest on their own as they establish certain regularities across projective transformations of the plane and introduce new view invariants (Lemma 4).

Theorem 2 is an obvious corollary of Theorem 1 but contains a significant practical aspect. It is shown that if the views ψ_1, ψ_2 are obtained by parallel projection, then F reduces to a special bilinear form — or, equivalently, that any perspective view ψ can be obtained by a rational linear function of two orthographic views. The reduction to a bilinear form implies that simpler recognition schemes are possible if the two reference views (model views) stored in memory are orthographic.

These two results may have several applications (discussed in Section 6), but the one emphasized throughout this paper is for the task of recognition of 3D objects via alignment. The alignment approach for recognition ([33, 16], and references therein) is based on the result that the equivalence class of views of an object (ignoring self occlusions) undergoing 3D rigid, affine or projective transformations can be captured by storing a 3D model of the object, or simply by storing at least two arbitrary “model” views of the object — assuming that the correspondence problem between the model views can somehow be solved (cf. [25, 5, 29]). During recognition a small number of corresponding points between the novel input view and the model views of a particular candidate object are sufficient to “re-project” the model onto the novel viewing position. Recognition is achieved if the re-projected image is successfully matched against the input image. We refer to the problem of predicting a novel view from a set of model views using a limited number of corresponding points, as the problem of *re-projection*.

The problem of re-projection can in principal be dealt with via 3D reconstruction of shape and camera motion. This includes classical structure from motion methods for recovering rigid camera motion parameters and metric shape [32, 18, 31, 14, 15], and more recent meth-

ods for recovering non-metric structure, i.e., assuming the objects undergo 3D affine or projective transformations, or equivalently, that the cameras are uncalibrated [17, 23, 35, 10, 13, 27, 28]. The classic approaches for perspective views are known to be unstable under errors in image measurements, narrow field of view, and internal camera calibration [3, 9, 12], and therefore, are unlikely to be of practical use for purposes of re-projection. The non-metric approaches, as a general concept, have not been fully tested on real images, but the methods proposed so far rely on recovering first the epipolar geometry — a process that is also known to be unstable in the presence of noise.

It is also known that the epipolar geometry is by itself sufficient to achieve re-projection by means of intersecting epipolar lines [22, 6, 8, 24, 21, 11]. This, however, is possible only if the centers of the three cameras are non-collinear — which can lead to numerical instability unless the centers are far from collinear — and any object point on the tri-focal plane cannot be re-projected as well. Furthermore, as with the non-metric reconstruction methods, obtaining the epipolar geometry is at best a sensitive process even when dozens of corresponding points are used and with the state of the art methods (see Section 5 for more details and comparative analysis with simulated and real images).

For purposes of stability, therefore, it is worthwhile exploring more direct tools for achieving re-projection. For instance, instead of reconstruction of shape and invariants we would like to establish a direct connection between views expressed as a functions of image coordinates alone — which we will call “algebraic functions of views”. Such a result was established in the orthographic case by [34]. There it was shown that any three orthographic views of an object satisfy a linear function of the corresponding image coordinates — this we will show here is simply a limiting case of larger set of algebraic functions, that in general have a trilinear form. With these functions one can manipulate views of an object, such as create new views, without the need to recover shape or camera geometry as an intermediate step — all what is needed is to appropriately combine the image coordinates of two reference views. Also, with these functions, the epipolar geometries are intertwined, leading not only to absence of singularities, but as we shall see in the experimental section to more accurate performance in the presence of errors in image measurements.

2 Notations

We consider object space to be the three-dimensional projective space \mathcal{P}^3 , and image space to be the two-dimensional projective space \mathcal{P}^2 . Let $\Phi \subset \mathcal{P}^3$ be a set of points standing for a 3D object, and let $\psi_i \subset \mathcal{P}^2$ denote views (arbitrary), indexed by i , of Φ . Given two cameras with centers located at $O, O' \in \mathcal{P}^3$, respectively, the epipoles are defined to be at the intersection of the line $\overline{OO'}$ with both image planes. Because the image plane is finite, we can assign, without loss of generality, the value 1 as the third homogeneous coordinate to every *observed* image point. That is, if (x, y) are the observed image co-

ordinates of some point (with respect to some arbitrary origin — say the geometric center of the image), then $p = (x, y, 1)$ denotes the homogeneous coordinates of the image plane. Since we will be working with at most three views at a time, we denote the relevant epipoles as follows: let $v \in \psi_1$ and $v' \in \psi_2$ be the corresponding epipoles between views ψ_1, ψ_2 , and let $\bar{v} \in \psi_1$ and $v'' \in \psi_3$ the corresponding epipoles between views ψ_1, ψ_3 . Likewise, corresponding image points across three views will be denoted by $p = (x, y, 1), p' = (x', y', 1)$ and $p'' = (x'', y'', 1)$. The term “image coordinates” will denote the non-homogeneous coordinate representation of \mathcal{P}^2 , e.g., $(x, y), (x', y'), (x'', y'')$ for the three corresponding points.

Planes will be denoted by π_i , indexed by i , and just π if only one plane is discussed. All planes are assumed to be arbitrary and distinct from one another. The symbol \cong denotes equality up to a scale, GL_n stands for the group of $n \times n$ matrices, and PGL_n is the group defined up to a scale.

A coordinate representation \mathcal{R} of \mathcal{P}^3 is a tetrad of coordinates $[z_o, z_1, z_2, z_3]$ such that if \mathcal{R}_0 is any one allowable representation, the whole class \mathcal{R} consists of all those representations that can be obtained from \mathcal{R}_0 by the action of the group PGL_4 . Given a set of views ψ_i , $i = 1, 2, \dots$, of Φ , where coordinates on ψ_1 are $[x, y, 1]$ and \mathcal{R}_0 is a representation for which $(z_o, z_1, z_2) = (x, y, 1)$, we will say that the object is undergoing at most 3D *relative affine transformations* between views if the class of representations \mathcal{R} consists of all those representations that can be obtained from \mathcal{R}_0 by the action of an *affine* subgroup of PGL_4 . In other words, the object undergoes some projective transformation and projected onto the view ψ_1 , after which all other transformations applied to Φ are affine. Note that this definition is general and allows full uncalibrated pin-hole camera motion (for more details on uncalibrated camera motion versus relative affine transformation versus taking pictures of pictures of the scene, see Appendix of [26]).

3 The Trilinear Form

The central result of this paper is presented in the following theorem. The remaining of the section is devoted to the proof of this result and its implications.

Theorem 1 (Trilinearity) *Let ψ_1, ψ_2, ψ_3 be three arbitrary perspective views of some object, modeled by a set of points in 3D, undergoing at most a 3D relative affine transformations between views. The image coordinates $(x, y) \in \psi_1, (x', y') \in \psi_2$ and $(x'', y'') \in \psi_3$ of three corresponding points across three views satisfy a pair of trilinear equations of the following form:*

$$\begin{aligned} x''(\alpha_1 x + \alpha_2 y + \alpha_3) + x'x''(\alpha_4 x + \alpha_5 y + \alpha_6) + \\ x'(\alpha_7 x + \alpha_8 y + \alpha_9) + \alpha_{10}x + \alpha_{11}y + \alpha_{12} = 0, \end{aligned}$$

and

$$\begin{aligned} y''(\beta_1 x + \beta_2 y + \beta_3) + y'y''(\beta_4 x + \beta_5 y + \beta_6) + \\ x'(\beta_7 x + \beta_8 y + \beta_9) + \beta_{10}x + \beta_{11}y + \beta_{12} = 0, \end{aligned}$$

where the coefficients $\alpha_j, \beta_j, j = 1, \dots, 12$, are fixed for all points, are uniquely defined up to an overall scale, and $\alpha_j = \beta_j, j = 1, \dots, 6$.

The following auxiliary propositions are used as part of the proof.

Lemma 1 (Auxiliary - Existence) *Let $A \in PGL_3$ be the projective mapping (homography) $\psi_1 \mapsto \psi_2$ due to some plane π . Let A be scaled to satisfy $p'_o \cong Ap_o + v'$, where $p_o \in \psi_1$ and $p'_o \in \psi_2$ are corresponding points coming from an arbitrary point $P_o \notin \pi$. Then, for any corresponding pair $p \in \psi_1$ and $p' \in \psi_2$ coming from an arbitrary point $P \in \mathcal{P}^3$, we have*

$$p' \cong Ap + kv'.$$

The coefficient k is independent of ψ_2 , i.e., is invariant to the choice of the second view.

The lemma, its proof and its theoretical and practical implications are discussed in detail in [26]. Note that the particular case where the homography A is affine, and the epipole v' is on the line at infinity, corresponds to the construction of affine structure from two orthographic views [17]. The scalar k is called a *relative affine invariant* and represents the ratio of the distance of P from π along the line of sight, and the distance of P from the camera center of ψ_1 , normalized by the ratio of distances of P_o from the plane and the camera center. This normalized ratio can be computed with the aid of a second arbitrary view ψ_2 .

Definition 1 *Homographies $A_i \in PGL_3$ from $\psi_1 \mapsto \psi_i$ due to the same plane π , are said to be scale-compatible if they are scaled to satisfy Lemma 1, i.e., for any point $P \in \Phi$ projecting onto $p \in \psi_1$ and $p^i \in \psi_i$, there exists a scalar k that satisfies*

$$p^i \cong A_i p + kv^i,$$

for any view ψ_i , where $v^i \in \psi_i$ is the epipole with ψ_1 (scaled arbitrarily).

Lemma 2 (Auxiliary — Uniqueness) *Let $A, A' \in PGL_3$ be two homographies of $\psi_1 \mapsto \psi_2$ due to planes π_1, π_2 , respectively. Then, there exists a scalar s , that satisfies the equation:*

$$A - sA' = [\alpha v', \beta v', \gamma v'],$$

for some coefficients α, β, γ .

Proof. Let $q \in \psi_1$ be any point in the first view. There exists a scalar s_q that satisfies $v' \cong Aq - s_q A'q$. Let $H = A - s_q A'$, and we have $Hq \cong v'$. But, as shown in [27], $Av \cong v'$ for any homography $\psi_1 \mapsto \psi_2$ due to any plane. Therefore, $Hv \cong v'$ as well. The mapping of two distinct points q, v onto the same point v' could happen only if $Hp \cong v'$ for all $p \in \psi_1$, and s_q is a fixed scalar s . This, in turn, implies that H is a matrix whose columns are multiples of v' . \square

Lemma 3 (Auxiliary for Lemma 4) *Let $A, A' \in PGL_3$ be homographies from $\psi_1 \mapsto \psi_2$ due to distinct planes π_1, π_2 , respectively, and $B, B' \in PGL_3$ be homographies from $\psi_1 \mapsto \psi_3$ due to π_1, π_2 , respectively. Then, $A' = AT$ for some $T \in PGL_3$, and $B = BCTC^{-1}$, where $Cv \cong \bar{v}$.*

Proof. Let $A = A_2^{-1}A_1$, where A_1, A_2 are homographies from ψ_1, ψ_2 onto π_1 , respectively. Similarly $B = B_2^{-1}B_1$, where B_1, B_2 are homographies from ψ_1, ψ_3 onto π_1 , respectively. Let $A_1\bar{v} = (c_1, c_2, c_3)^T$, and let $C \cong A_1^{-1}diag(c_1, c_2, c_3)A_1$. Then, $B_1 \cong A_1C^{-1}$, and thus, we have $B \cong B_2^{-1}A_1C^{-1}$. Note that the only difference between A_1 and B_1 is due to the different location of the epipoles v, \bar{v} , which is compensated by C ($Cv \cong \bar{v}$). Let $E_1 \in PGL_3$ be the homography from ψ_1 to π_2 , and $E_2 \in PGL_3$ the homography from π_2 to π_1 . Then with proper scaling of E_1 and E_2 we have

$$A' = A_2^{-1}E_2E_1 = AA_1^{-1}E_2E_1 = AT,$$

and with proper scaling of C we have,

$$B' = B_2^{-1}E_2E_1C^{-1} = BCA_1^{-1}E_2E_1C^{-1} = BCTC^{-1}.$$

□

Lemma 4 (Auxiliary — Uniqueness)

For scale-compatible homographies, the scalars s, α, β, γ of Lemma 2 are invariants indexed by ψ_1, π_1, π_2 . That is, given an arbitrary third view ψ_3 , let B, B' be the homographies from $\psi_1 \mapsto \psi_3$ due to π_1, π_2 , respectively. Let B be scale-compatible with A , and B' be scale-compatible with A' . Then,

$$B - sB' = [\alpha v'', \beta v'', \gamma v''].$$

Proof. We show first that s is invariant, i.e., that $B - sB'$ is a matrix whose columns are multiples of v'' . From Lemma 2, and Lemma 3 there exists a matrix H , whose columns are multiples of v' , a matrix T that satisfies $A' = AT$, and a scalar s such that $I - sT = A^{-1}H$. After multiplying both sides by BC , and then pre-multiplying by C^{-1} we obtain

$$B - sBCTC^{-1} = BCA^{-1}HC^{-1}.$$

From Lemma 3, we have $B' = BCTC^{-1}$. The matrix $A^{-1}H$ has columns which are multiples of v (because $A^{-1}v' \cong v$), $CA^{-1}H$ is a matrix whose columns are multiple of \bar{v} , and $BCA^{-1}H$ is a matrix whose columns are multiples of v'' . Pre-multiplying $BCA^{-1}H$ by C^{-1} does not change its form because every column of $BCA^{-1}HC^{-1}$ is simply a linear combination of the columns of $BCA^{-1}H$. As a result, $B - sB'$ is a matrix whose columns are multiples of v'' .

Let $H = A - sA'$ and $\hat{H} = B - sB'$. Since the homographies are scale compatible, we have from Lemma 1 the existence of invariants k, k' associated with an arbitrary $p \in \psi_1$, where k is due to π_1 , and k' is due to π_2 : $p' \cong Ap + kv' \cong A'p + k'v'$ and $p'' \cong Bp + kv'' \cong B'p + k'v''$. Then from Lemma 2 we have $Hp = (sk' - k)v'$ and $\hat{H}p = (sk' - k)v''$. Since p is arbitrary, this could happen only if the coefficients of the multiples of v' in H and the coefficients of the multiples of v'' in \hat{H} , coincide.

□

Proof of Theorem: Lemma 1 provides the existence part of theorem, as follows. Since Lemma 1 holds for any plane, choose a plane π_1 and let A, B be the scale-compatible homographies $\psi_1 \mapsto \psi_2$ and $\psi_1 \mapsto \psi_3$, respectively. Then, for every point $p \in \psi_1$, with corresponding

points $p' \in \psi_2, p'' \in \psi_3$, there exists a scalar k that satisfies: $p' \cong Ap + kv'$, and $p'' \cong Bp + kv''$. We can isolate k from both equations and obtain:

$$k = \frac{v'_1 - x'v'_3}{(x'a_3 - a_1)^T p} = \frac{v'_2 - y'v'_3}{(y'a_3 - a_2)^T p} = \frac{y'v'_1 - x'v'_2}{(x'a_2 - y'a_1)^T p}, \quad (1)$$

$$k = \frac{v''_1 - x''v''_3}{(x''b_3 - b_1)^T p} = \frac{v''_2 - y''v''_3}{(y''b_3 - b_2)^T p} = \frac{y''v''_1 - x''v''_2}{(x''b_2 - y''b_1)^T p}, \quad (2)$$

where $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3$ and $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$ are the row vectors of A and B and $v' = (v'_1, v'_2, v'_3)$, $v'' = (v''_1, v''_2, v''_3)$. Because of the invariance of k we can equate terms of Equation 1 with terms of Equation 2 and obtain trilinear functions of image coordinates across three views. For example, by equating the first two terms in each of the equations, we obtain:

$$\begin{aligned} x''(v'_1\mathbf{b}_3 - v'_3\mathbf{a}_1)^T p + x''x'(v'_3\mathbf{a}_3 - v'_3\mathbf{b}_3)^T p + \\ x'(v'_3\mathbf{b}_1 - v'_1\mathbf{a}_3)^T p + (v'_1\mathbf{a}_1 - v'_1\mathbf{b}_1)^T p = 0, \end{aligned} \quad (3)$$

In a similar fashion, after equating the first term of Equation 1 with the second term of Equation 2, we obtain:

$$\begin{aligned} y''(v'_1\mathbf{b}_3 - v'_3\mathbf{a}_1)^T p + y''x'(v'_3\mathbf{a}_3 - v'_3\mathbf{b}_3)^T p + \\ x'(v'_3\mathbf{b}_2 - v'_2\mathbf{a}_3)^T p + (v'_2\mathbf{a}_1 - v'_1\mathbf{b}_2)^T p = 0. \end{aligned} \quad (4)$$

Both equations are of the desired form, with the first six coefficients identical across both equations.

The question of uniqueness arises because Lemma 1 holds for any plane. If we choose a different plane, say π_2 , with homographies A', B' , then we must show that the new homographies give rise to the same coefficients (up to an overall scale). The parenthesized terms in Equations 3 and 4 have the general form: $v'_i\mathbf{b}_i \pm v'_j\mathbf{a}_j$, for some i and j . Thus, we need to show that there exists a scalar s that satisfies

$$v''_i(\mathbf{a}_j - s\mathbf{a}'_j) = v'_j(\mathbf{b}_i - s\mathbf{b}'_i).$$

This, however, follows directly from Lemmas 2 and 4. □

The direct implication of the theorem is that one can generate a novel view (ψ_3) by simply combining two model views (ψ_1, ψ_2). The coefficients α_j and β_j of the combination can be recovered together as a solution of a linear system of 17 equations ($24 - 6 - 1$) given nine corresponding points across the three views (more than nine points can be used for a least-squares solution).

Taken together, Equations 1 and 2 lead to 9 algebraic functions of three views, six of which are separate for x'' and y'' . The other four functions are listed below:

$$\begin{aligned} x''(\cdot) + x''y'(\cdot) + y'(\cdot) + (\cdot) &= 0, \\ y''(\cdot) + y''y'(\cdot) + y'(\cdot) + (\cdot) &= 0, \\ x''x'(\cdot) + x''y'(\cdot) + x'(\cdot) + y'(\cdot) &= 0, \\ y''x'(\cdot) + y''y'(\cdot) + x'(\cdot) + y'(\cdot) &= 0, \end{aligned}$$

where (\cdot) represent linear polynomials in x, y . The solution for x'', y'' is unique without constraints on the allowed camera transformations. If we choose Equations 3 and 4, then v'_1 and v'_3 should not vanish simultaneously, i.e., $v' \cong (0, 1, 0)$ is a singular case. Also $v'' \cong (0, 1, 0)$ and $v'' \cong (1, 0, 0)$ give rise to singular cases. One can easily show that for each singular case there are two other functions out of the nine available ones that provide a

unique solution for x'', y'' . Note that the singular cases are pointwise, i.e., only three epipolar directions are excluded, compared to the more wide-spread singular cases that occur with epipolar intersection, as described in the introduction.

In practical terms, the process of generating a novel view can be easily accomplished without the need to explicitly recover structure, camera transformation, or just the epipolar geometry. The process described here is fundamentally different from intersecting epipolar lines in the following ways: first, we use the three views together, instead of pairs of views separately; second, there is no process of line intersection, i.e., the x and y coordinates of ψ_3 are obtained separately as a solution of a single equation in coordinates of the other two views; and thirdly, the process is well defined in cases where intersecting epipolar lines becomes singular (e.g., when the three camera centers are collinear). Furthermore, by avoiding the need to recover the epipolar geometry we obtain a significant practical advantage, since the epipolar geometry is the most error-sensitive component when working with perspective views.

The connection between the general result of trilinear functions of views to the “linear combination of views” result [34] for orthographic views, can easily be seen by setting A and B to be affine in \mathcal{P}^2 , and $v'_3 = v''_3 = 0$. For example, Equation 3 reduces to

$$v'_1 x'' - v''_1 x' + (v''_1 \mathbf{a}_1 - v'_1 \mathbf{b}_1)^T p = 0,$$

which is of the form

$$\alpha_1 x'' + \alpha_2 x' + \alpha_3 x + \alpha_4 y + \alpha_5 = 0.$$

Thus, in the case where all three views are orthographic, x'' is expressed as a linear combination of image coordinates of the two other views — as discovered by [34].

4 The Bilinear Form

Consider the case for which the two reference (model) views of an object are taken orthographically (using a tele lens would provide a reasonable approximation), but during recognition any perspective view of the object is allowed. It can easily be shown that the three views are then connected via bilinear functions (instead of trilinear):

Theorem 2 (Bilinearity) *Within the conditions of Theorem 1, in case the views ψ_1 and ψ_2 are obtained by parallel projection, then the pair of trilinear forms of Theorem 1 reduce to the following pair of bilinear equations:*

$$x''(\alpha_1 x + \alpha_2 y + \alpha_3) + \alpha_4 x'' x' + \alpha_5 x' + \alpha_6 x + \alpha_7 y + \alpha_8 = 0,$$

and

$$y''(\beta_1 x + \beta_2 y + \beta_3) + \beta_4 y'' x' + \beta_5 x' + \beta_6 x + \beta_7 y + \beta_8 = 0,$$

where $\alpha_j = \beta_j$, $j = 1, \dots, 4$.

Proof. Under these conditions we have from Lemma 1 that A is affine in \mathcal{P}^2 and $v'_3 = 0$, therefore Equation 3 reduces to:

$$x''(v'_1 \mathbf{b}_3 - v''_3 \mathbf{a}_1)^T p + v''_3 x'' x' - v''_1 x' + (v''_1 \mathbf{a}_1 - v'_1 \mathbf{b}_1)^T p = 0.$$

Similarly, Equation 4 reduces to:

$$y''(v'_1 \mathbf{b}_3 - v''_3 \mathbf{a}_1)^T p + v''_3 y'' x' - v''_2 x' + (v''_2 \mathbf{a}_1 - v'_1 \mathbf{b}_2)^T p = 0.$$

Both equations are of the desired form, with the first four coefficients identical across both equations. \square

A bilinear function of three views has two advantages over the general trilinear function. First, only six corresponding points (instead of nine) across three views are required for solving for the coefficients. Second, the lower the degree of the algebraic function, the less sensitive the solution should be in the presence of errors in measuring correspondences. In other words, it is likely (though not necessary) that the higher order terms, such as the term $x'' x' x$ in Equation 3, will have a higher contribution to the overall error sensitivity of the system.

Compared to the case when all views are assumed orthographic, this case is much less of an approximation. Since the model views are taken only once, it is not unreasonable to require that they be taken in a special way, namely, with a tele lens (assuming we are dealing with object recognition, rather than scene recognition). If that requirement is satisfied, then the recognition task is general since we allow any perspective view to be taken during the recognition process.

5 Experimental Data

The experiments described in this section were done in order to evaluate the practical aspect of using the trilinear result for re-projection compared to using epipolar intersection and the linear combination result of [34] (the latter we have shown is simply a limiting case of the trilinear result).

The epipolar intersection method was implemented in the following way. Let F_{13} and F_{23} be the matrices (“essential” matrices in classical terminology [18], which we adopt here) that satisfy $p'' F_{13} p = 0$, and $p'' F_{23} p' = 0$. Then, by incidence of p'' with its epipolar line, we have:

$$p'' \cong F_{13} p \times F_{23} p'.$$

Therefore, given eight corresponding points across the three views, we can recover the two essential matrices, and then re-project all other object points onto the third view. In practice one would use more than eight points for recovering the essential matrices in a linear or non-linear squares method. Since linear least squares methods are still sensitive to image noise, we used the implementation of a non-linear method described in [19] which was kindly provided by T. Luong and L. Quan.

The first experiment is with simulation data showing that even when the epipolar geometry is recovered accurately, it is still significantly better to use the trilinear result which avoids the process of line intersection. The second experiment is done on a real set of images, comparing the performance of the various methods and the number of corresponding points that are needed in practice to achieve reasonable re-projection results.

5.1 Computer Simulations

We used an object of 46 points placed randomly with z coordinates between 100 units and 120 units, and x, y

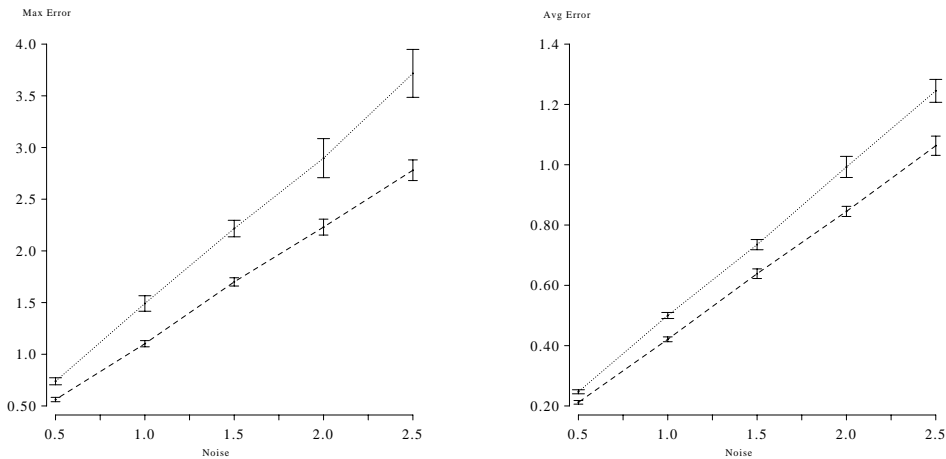


Figure 1: Comparing the performance of the epipolar intersection method (the dotted line) and the trilinear functions method (dashed line) in the presence of image noise. The graph on the left shows the maximal re-projection error averaged over 200 trials per noise level (bars represent standard deviation). Graph on the right displays the average re-projection error averaged over all re-projected points averaged over the 200 trials per noise level.

coordinates ranging randomly between -125 and $+125$. Focal length was of 50 units and the first view was obtained by $f_x/z, f_y/z$. The second view (ψ_2) was generated by a rotation around the point $(0, 0, 100)$ with axis $(0.14, 0.7, 0.7)$ and by an angle of 0.3 radians. The third view (ψ_3) was generated by a rotation around an axis $(0, 1, 0)$ with the same translation and angle. Various amounts of random noise was applied to all points that were to be re-projected onto a third view, but not to the eight or nine points that were used for recovering the parameters (essential matrices, or trilinear coefficients). The noise was random, added separately to each coordinate and with varying levels from 0.5 to 2.5 pixel error. We have done 1000 trials as follows: 20 random objects were created, and for each degree of error the simulation was ran 10 times per object. We collected the maximal re-projection error (in pixels) and the average re-projection error (averaged of all the points that were re-projected). These numbers were collected separately for each degree of error by averaging over all trials (200 of them) and recording the standard deviation as well. Since no error were added to the eight or nine points that were used to determine the epipolar geometry and the trilinear coefficients, we simply solved the associated linear systems of equations required to obtain the essential matrices or the trilinear coefficients.

The results are shown in Figure 1. The graph on the left shows the performance of both algorithms for each level of image noise by measuring the maximal re-projection error. We see that under all noise levels, the trilinear method is significantly better and also has a smaller standard deviation. Similarly for the average re-projection error shown in the graph on the right.

This difference in performance is expected, as the trilinear method takes all three views together, rather than every pair separately, and thus avoiding line intersections.

5.2 Experiments On Real Images

Figure 2 shows three views of the object we selected for the experiment. The object is a sports shoe with added texture to facilitate the correspondence process. This object was chosen because of its complexity, i.e., it has a shape of a natural object and cannot easily be described parametrically (as a collection of planes or algebraic surfaces). Note that the situation depicted here is challenging because the re-projected view is not in-between the two model views, i.e., one should expect a larger sensitivity to image noise than in-between situations. A set of 34 points were manually selected on one of the frames, ψ_1 , and their correspondences were automatically obtained along all other frames used in this experiment. The correspondence process is based on an implementation of a coarse-to-fine optical-flow algorithm described in [7]. To achieve accurate correspondences across distant views, intermediate in-between frames were taken and the displacements across consecutive frames were added. The overall displacement field was then used to push (“warp”) the first frame towards the target frame and thus create a synthetic image. Optical-flow was applied again between the synthetic frame and the target frame and the resulting displacement was added to the overall displacement obtained earlier. This process provides a dense displacement field which is then sampled to obtain the correspondences of the 34 points initially chosen in the first frame. The results of this process are shown in Figure 2 by displaying squares centered around the computed locations of the corresponding points. One can see that the correspondences obtained in this manner are reasonable, and in most cases to sub-pixel accuracy. One can readily automate further this process by selecting points in the first frame for which the Hessian matrix of spatial derivatives is well conditioned — similar to the confidence values suggested in the implementations of [4, 7, 30] — however, the intention here was not so much as to build a complete system but to test the

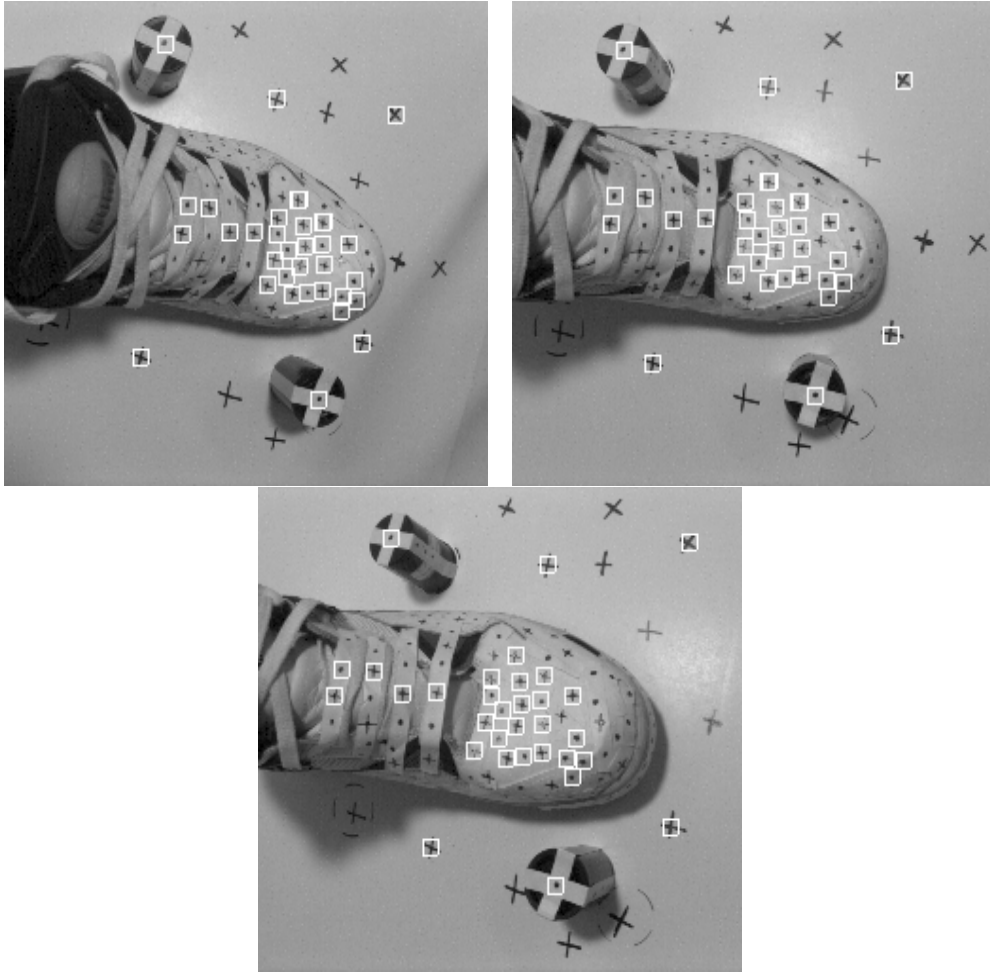


Figure 2: *Top Row*: Two model views, ψ_1 on the left and ψ_2 on the right. The overlaid squares illustrate the corresponding points (34 points). *Bottom Row*: Third view ψ_3 . Note that ψ_3 is not in-between ψ_1 and ψ_2 , making the re-projection problem more challenging (i.e., performance is more sensitive to image noise than in-between situations).

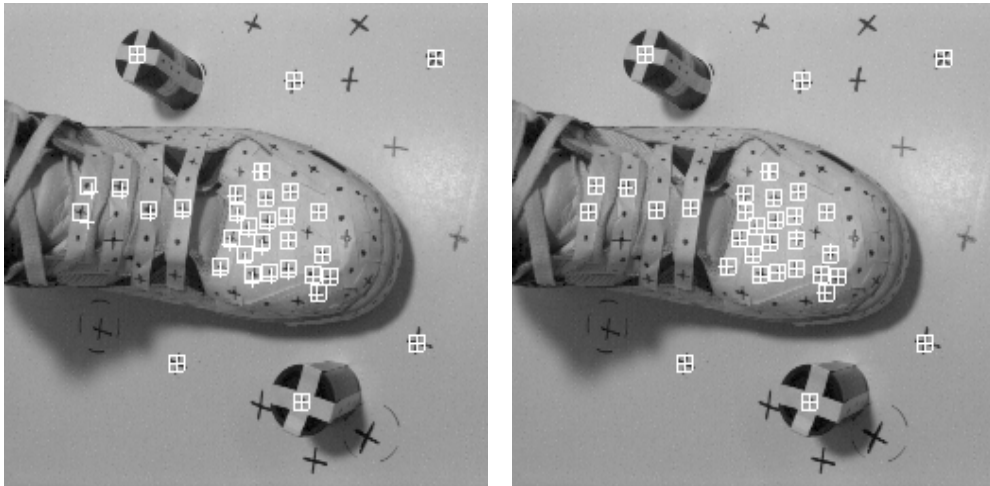


Figure 3: Re-projection onto ψ_3 using the trilinear result. The re-projected points are marked as crosses, therefore should be at the center of the squares for accurate re-projection. On the left, the minimal number of points were used for recovering the trilinear coefficients (nine points); the average pixel error between the true and estimated locations is 1.4, and the maximal error is 5.7. On the right 12 points were used in a least squares fit; average error is 0.4 and maximal error is 1.4.

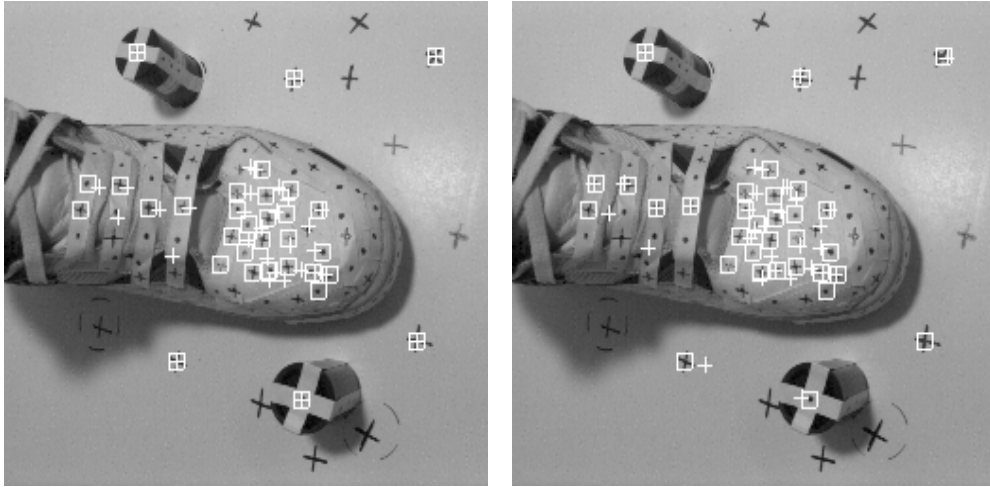


Figure 4: Results of re-projection using intersection of epipolar lines. The re-projected points are marked as crosses, therefore should be at the center of the squares for accurate re-projection. In the lefthand display the ground plane points were used for recovering the essential matrix (see text), and in the righthand display the essential matrices were recovered from the implementation of [19] using all 34 points across the three views. Maximum displacement error in the lefthand display is 25.7 pixels and average error is 7.7 pixels. Maximal error in the righthand display is 43.4 pixels and average error is 9.58 pixels.

performance of the trilinear re-projection method and compare it to the performance of epipolar intersection and the linear combination methods.

The trilinear method requires at least nine corresponding points across the three views (we need 17 equation, and nine points provide 18 equations), whereas epipolar intersection can be done (in principle) with eight points. The question we are about to address is what is the number of points that are required in practice (due to errors in correspondence, lens distortions and other effects that are not adequately modeled by the pin-hole camera model) to achieve reasonable performance?

The trilinear result was first applied with the minimal number of points (nine) for solving for the coefficients, and then applied with 12 points using a linear least-squares solution. The results are shown in Figure 3. Nine points provide a re-projection with maximal error of 5.7 pixels and average error of 1.4 pixels. The solution using 12 points provided a significant improvement with maximal error of 1.4 and average error of 0.4 pixels. Using more points did not improve significantly the results; for example, when all 34 points were used the maximal error went down to 1.14 pixels and average error stayed at 0.42 pixels.

Next the epipolar intersection method was applied. We used two methods for recovering the essential matrices. One method is by using the implementation of [19], and the other is by taking advantage that four of the corresponding points are coming from a plane (the ground plane). In the former case, much more than eight points were required in order to achieve reasonable results. For example, when using all the 34 points, the maximal error was 43.4 pixels and the average error was 9.58 pixels. In the latter case, we recovered first the homography B due to the ground plane and then the epipole v'' using two additional points (those on the film cartridges). It

is then known (see [26, 20]) that $F_{13} = [v'']B$, where $[v'']$ is the anti-symmetric matrix of v'' . A similar procedure was used to recover F_{23} . Therefore, only six points were used for re-projection, but nevertheless, the results were slightly better: maximal error of 25.7 pixels and average error of 7.7 pixels. Figure 4 shows these results.

Finally, we tested the performance of re-projection using the linear combination method. Since the linear combination methods holds only for orthographic views, we are actually testing the orthographic assumption under a perspective situation, or in other words, whether the higher (bilinear and trilinear) order terms of the trilinear equations are significant or not. The linear combination method requires at least four corresponding points across the three views. We applied the method with four, 12 (for comparison with the trilinear case shown in Figure 3), and all 34 points (the latter two using linear least squares). The results are displayed in Figure 5. The performance in all cases are significantly poorer than when using the trilinear functions, but better than the epipolar intersection method.

6 Discussion

We have seen that any view of a fixed 3D object can be expressed as a trilinear function with two reference views in the general case, or as a bilinear function when the reference views are created by means of parallel projection. These functions provide alternative, much simpler, means for manipulating views of a scene than other methods. Experimental results show that the trilinear functions are also useful in practice yielding performance that is significantly better than epipolar intersection or the linear combination method.

The application that was emphasized throughout the paper is visual recognition via alignment. Reasonable

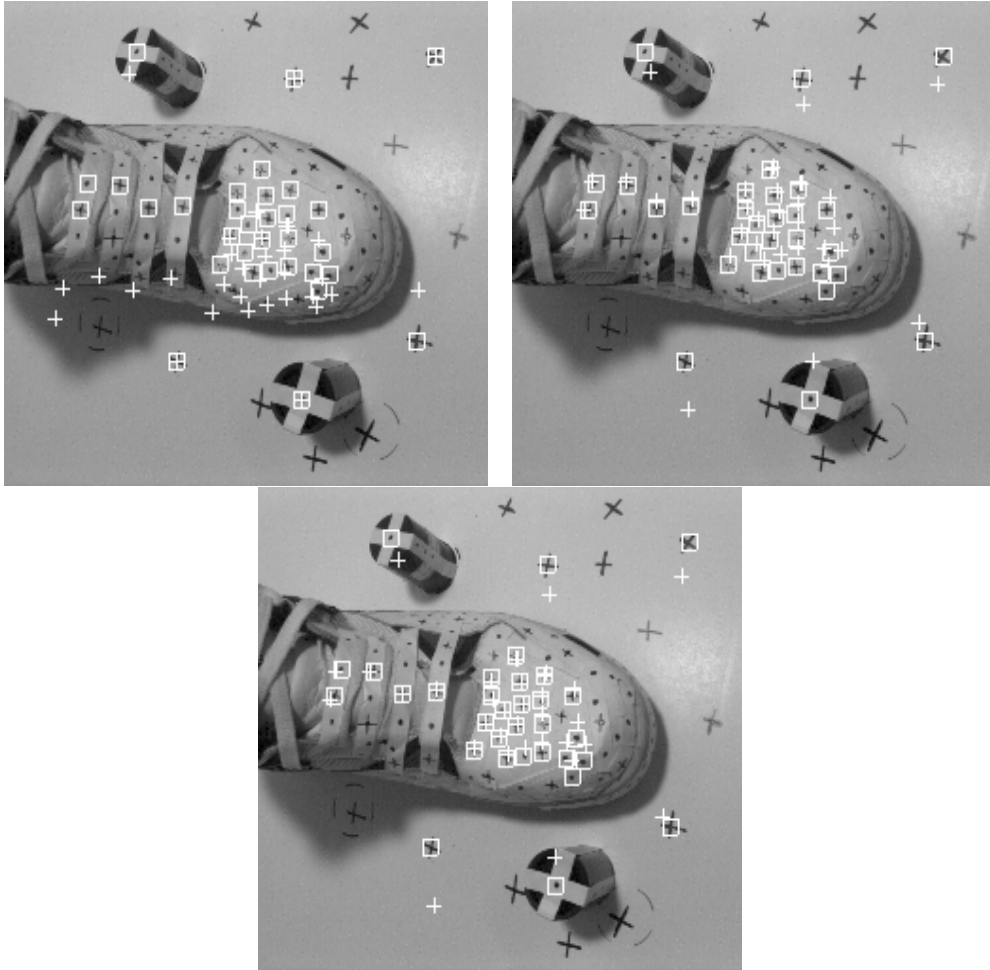


Figure 5: Results of re-projection using the linear combination of views method proposed by [34] (applicable to parallel projection). *Top Row:* In the lefthand display the linear coefficients were recovered from four corresponding points; maximal error is 56.7 pixels and average error is 20.3 pixels. In the righthand display the coefficients were recovered using 12 points in a linear least squares fashion; maximal error is 24.3 pixels and average error is 6.8 pixels. *Bottom Row:* The coefficients were recovered using all 34 points across the three views. Maximal error is 29.4 pixels and average error is 5.03 pixels.

performance was obtained with 12 corresponding points with the novel view (ψ_3) — which may be too many if the image to model matching is done by trying all possible combinations of point matches. The existence of bilinear functions in the special case where the model is orthographic, but the novel view is perspective, is more encouraging from the standpoint of counting points. Here we have the result that only six corresponding points are required to obtain recognition of perspective views (provided we can satisfy the requirement that the model is orthographic). We have not experimented with bilinear functions to see how many points would be needed in practice, but plan to do that in the future. Because of their simplicity, one may speculate that these algebraic functions will find uses in tasks other than visual recognition — some of those are discussed below.

There may exist other applications where simplicity is of major importance, whereas the number of points is less of a concern. Consider for example, the application of model-based compression. With the trilinear functions we need 17 parameters to represent a view as a function of two reference views in full correspondence. Assume both the sender and the receiver have the two reference views and apply the same algorithm for obtaining correspondences between the two views. To send a third view (ignoring problems of self occlusions that could be dealt separately) the sender can solve for the 17 parameters using many points, but eventually send only the 17 parameters. The receiver then simply combines the two reference views in a “trilinear way” given the received parameters. This is clearly a domain where the number of points are not a major concern, whereas simplicity, and robustness (as shown above) due to the short-cut in the computations, is of great importance.

Related to image coding, an approach of image decomposition into “layers” was recently proposed by [1, 2]. In this approach, a sequence of views is divided up into regions, whose motion of each is described approximately by a 2D affine transformation. The sender sends the first image followed only by the six affine parameters for each region for each subsequent frame. The use of algebraic functions of views can potentially make this approach more powerful because instead of dividing up the scene into planes (it would have been planes if the projection was parallel, in general its not even planes) one can attempt to divide the scene into objects, each carries the 17 parameters describing its displacement onto the subsequent frame.

Another area of application may be in computer graphics. Re-projection techniques provide a short-cut for image rendering. Given two fully rendered views of some 3D object, other views (again ignoring self-occlusions) can be rendered by simply “combining” the reference views. Again, the number of corresponding points is less of a concern here.

Acknowledgments

Thanks to W.E.L. Grimson for critical reading of the report. I thank T. Luong and L. Quan for providing the implementation for recovering essential matrices and epipoles. Thanks to N. Navab and A. Azarbayejani for

assistance in capturing the image sequence (equipment courtesy of MIT Media Laboratory).

References

- [1] E.H. Adelson. Layered representations for image coding. Technical Report 181, Media Laboratory, Massachusetts Institute of Technology, 1991.
- [2] E.H. Adelson and J.Y.A. Wang. Layered representation for motion analysis. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition*, pages 361–366, New York, NY, June 1993.
- [3] G. Adiv. Inherent ambiguities in recovering 3D motion and structure from a noisy flow field. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-11(5):477–489, 1989.
- [4] P. Anandan. A unified perspective on computational techniques for the measurement of visual motion. In *Proceedings Image Understanding Workshop*, pages 219–230, Los Angeles, CA, February 1987. Morgan Kaufmann, San Mateo, CA.
- [5] I.A. Bachelder and S. Ullman. Contour matching using local affine transformations. In *Proceedings Image Understanding Workshop*. Morgan Kaufmann, San Mateo, CA, 1992.
- [6] E.B. Barrett, M.H. Brill, N.N. Haag, and P.M. Payton. Invariant linear methods in photogrammetry and model-matching. In J.L. Mundy and A. Zisserman, editors, *Applications of invariances in computer vision*. MIT Press, 1992.
- [7] J.R. Bergen and R. Hingorani. Hierarchical motion-based frame rate conversion. Technical report, David Sarnoff Research Center, 1990.
- [8] S. Demey, A. Zisserman, and P. Beardsley. Affine and projective structure from motion. In *Proceedings of the British Machine Vision Conference*, October 1992.
- [9] R. Dutta and M.A. Synder. Robustness of correspondence based structure from motion. In *Proceedings of the International Conference on Computer Vision*, pages 106–110, Osaka, Japan, December 1990.
- [10] O.D. Faugeras. What can be seen in three dimensions with an uncalibrated stereo rig? In *Proceedings of the European Conference on Computer Vision*, pages 563–578, Santa Margherita Ligure, Italy, June 1992.
- [11] O.D. Faugeras and L. Robert. What can two images tell us about a third one? Technical Report INRIA, France, 1993.
- [12] W.E.L. Grimson. Why stereo vision is not always about 3D reconstruction. A.I. Memo No. 1435, Artificial Intelligence Laboratory, Massachusetts Institute of Technology, July 1993.
- [13] R. Hartley, R. Gupta, and T. Chang. Stereo from uncalibrated cameras. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition*, pages 761–764, Champaign, IL., June 1992.

- [14] B.K.P. Horn. Relative orientation. *International Journal of Computer Vision*, 4:59–78, 1990.
- [15] B.K.P. Horn. Relative orientation revisited. *Journal of the Optical Society of America*, 8:1630–1638, 1991.
- [16] D.P. Huttenlocher and S. Ullman. Recognizing solid objects by alignment with an image. *International Journal of Computer Vision*, 5(2):195–212, 1990.
- [17] J.J. Koenderink and A.J. Van Doorn. Affine structure from motion. *Journal of the Optical Society of America*, 8:377–385, 1991.
- [18] H.C. Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293:133–135, 1981.
- [19] Q.T. Luong, R. Deriche, O.D. Faugeras, and T. Papadopoulos. On determining the fundamental matrix: Analysis of different methods and experimental results. Technical Report INRIA, France, 1993.
- [20] Q.T. Luong and T. Vieville. Canonical representations for the geometries of multiple projective views. Technical Report INRIA, France, 1993.
- [21] J. Mundy and A. Zisserman. Appendix — projective geometry for machine vision. In J. Mundy and A. Zisserman, editors, *Geometric invariances in computer vision*. MIT Press, Cambridge, 1992.
- [22] J.L. Mundy, R.P. Welty, M.H. Brill, P.M. Payton, and E.B. Barrett. 3-D model alignment without computing pose. In *Proceedings Image Understanding Workshop*, pages 727–735. Morgan Kaufmann, San Mateo, CA, January 1992.
- [23] A. Shashua. Correspondence and affine shape from two orthographic views: Motion and Recognition. A.I. Memo No. 1327, Artificial Intelligence Laboratory, Massachusetts Institute of Technology, December 1991.
- [24] A. Shashua. *Geometry and Photometry in 3D visual recognition*. PhD thesis, M.I.T Artificial Intelligence Laboratory, AI-TR-1401, November 1992.
- [25] A. Shashua. Illumination and view position in 3D visual recognition. In S.J. Hanson J.E. Moody and R.P. Lippmann, editors, *Advances in Neural Information Processing Systems 4*, pages 404–411. San Mateo, CA: Morgan Kaufmann Publishers, 1992. Proceedings of the fourth annual conference NIPS, Dec. 1991, Denver, CO.
- [26] A. Shashua. On geometric and algebraic aspects of 3D affine and projective structures from perspective 2D views. In *The 2nd European Workshop on Invariants*, Azores Islands, Portugal, October 1993. Also in MIT AI memo No. 1405, July 1993.
- [27] A. Shashua. Projective depth: A geometric invariant for 3D reconstruction from two perspective/orthographic views and for visual recognition. In *Proceedings of the International Conference on Computer Vision*, pages 583–590, Berlin, Germany, May 1993.
- [28] A. Shashua. Projective structure from uncalibrated images: structure from motion and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1994. in press.
- [29] A. Shashua and S. Toelg. The quadric reference surface: Applications in registering views of complex 3d objects. In *Proceedings of the European Conference on Computer Vision*, Stockholm, Sweden, May 1994.
- [30] C. Tomasi and T. Kanade. Factoring image sequences into shape and motion. In *IEEE Workshop on Visual Motion*, pages 21–29, Princeton, NJ, September 1991.
- [31] R.Y. Tsai and T.S. Huang. Uniqueness and estimation of three-dimensional motion parameters of rigid objects with curved surface. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-6:13–26, 1984.
- [32] S. Ullman. *The Interpretation of Visual Motion*. MIT Press, Cambridge and London, 1979.
- [33] S. Ullman. Aligning pictorial descriptions: an approach to object recognition. *Cognition*, 32:193–254, 1989. Also: in MIT AI Memo 931, Dec. 1986.
- [34] S. Ullman and R. Basri. Recognition by linear combination of models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-13:992–1006, 1991. Also in M.I.T AI Memo 1052, 1989.
- [35] D. Weinshall. Model based invariants for 3-D vision. *International Journal of Computer Vision*, 10(1):27–42, 1993.