

# Placing Three-Dimensional Models in an Uncalibrated Single Image of an Architectural Scene

Sara Keren

Electrical Engineering Dept.  
Technion

Ilan Shimshoni\*

Industrial Engineering Dept.  
Technion  
ilans@ie.technion.ac.il

Ayellet Tal<sup>†</sup>

Electrical Engineering Dept.  
Technion  
ayellet@ee.technion.ac.il

Corresponding author: Ilan Shimshoni, phone: +972-4-8294415, fax: +972-4-8294688

**Keywords:** Augmented reality, camera calibration, architectural scenes.

---

\*Ilan Shimshoni was supported in part by The Israeli Ministry of Science grant no. 2104

<sup>†</sup>Ayellet Tal was supported in part by Israeli Ministry of Science grant no. 01-01-01509.

This paper discusses the problem of inserting three-dimensional models into a single image. The main focus of the paper is on the accurate recovery of the camera’s parameters, so that 3D models can be inserted in the “correct” position and orientation. The paper addresses two issues. The first is an automatic extraction of the principal vanishing points from an image. The second is a theoretical and an experimental analysis of the errors. To test the concept, a system which “plants” virtual 3D objects in the image was implemented. It was tested on many indoor augmented reality scenes. Our analysis and experiments have shown that errors in the placement of the objects are un-noticeable.

## 1 Introduction

A basic problem in many augmented reality systems is how to create new artificial images given real images [Azuma, 1997]. A common approach is to create a new image of a scene from several other images of that scene taken from different viewpoints [Chang et al., 1999, Shade et al., 1998]. Another approach is to insert 3D models into video sequences [Berger et al., 1999, Mellor, 1995, de Agapito et al., 1999, Kutulakos and Vallino, 1998]. The approach we are pursuing in this paper is that of inserting 3D models into a single image.

Solving this problem requires solving several sub-problems. The first sub-problem is how to recover the camera’s internal calibration parameters and position within the scene accurately enough so that 3D models can be inserted into the image in the “correct” position [Hartley and Zisserman, 2000, Caprile and Torre, 1990, Criminisi et al., 2000]. A standard way to handle it, used in several augmented reality systems, is to place a known easily-detectable 3D object in the scene, and to recover the projection matrix from it [Kato and Billinghurst, 1999, MacIntyre and Coelho, 2002, Poupyrev et al., 2002]. We will show how to do it without adding calibration objects into the scene. The second sub-problem is how to recover the lighting conditions in the scene [Boivin and Gagalowicz, 2001, Koudelka et al., 2001, Debevec, 1998, Ramamoorthi and Hanrahan, 2001]. The third sub-problem is how to deal correctly with occlusions of the new object with the existing objects in the scene [Wloka and Anderson, 1995].

In this paper we concentrate on the first sub-problem. Figure 1 illustrates the intention of our algorithm and the results of running our system on a typical image. The original picture was taken by a digital camera whose internal parameters and its location and orientation in the room are unknown. We “placed” on the table artificial objects. The green chime

“hanging” from the ceiling is also artificial. We have chosen to place objects which look artificial, to distinguish them from real objects. The emphasis here is on the perfect location, orientation and size of the objects in the image.



Figure 1: An augmented room

Solving the problem of placing artificial objects given a single image without any prior knowledge about the scene is impossible. We therefore chose a large class of images for which there is prior information about the class but not much information has to be supplied with the image in order to perform the task. This class contains images of architectural scenes [Criminisi et al., 2000, Liebowitz et al., 1999, Antone and Teller, 2001, Teller et al., 2001] and in particular the interior of buildings. Typically, in these scenes lines appearing in the images are parallel to the three principal axes. This information can be used to recover the required parameters.

Still the problem is not trivial. The main question is the following. Can the unknown parameters be recovered accurately enough from noisy measurements from one image such that they can be used to “plant” 3D models in the image in a way that they will look correct to the viewer. In this paper we answer this question in the affirmative.

Basically, the problem is how to estimate from the image the projection matrix which transforms 3D points to points in the image. In most augmented reality systems this matrix is estimated from the correspondence between known 3D points and points detected in the

image [MacIntyre and Coelho, 2002]. Estimating this projection matrix from lines parallel to the principal axes was studied in [Hartley and Zisserman, 2000, Caprile and Torre, 1990, Criminisi et al., 2000]. Usually to evaluate such algorithms, the accuracy of the recovered projection matrix is compared to the ground-truth [MacIntyre and Coelho, 2002]. However, in our application the most important question is whether the results are accurate enough so that a human viewer will not detect differences between a real object and its 3D model inserted into the scene. We therefore direct our theoretical and experimental analysis to test projection errors.

There are two possible ways to insert 3D models into images. The first option is to reconstruct a three-dimensional scene from the images and let the user “place” the artificial objects in their locations in three dimensions. Then the augmented scene is projected back to image space. The second option is to choose a location in the image for the object to be placed in. Consequentially, the errors measuring the distance between the real location and the calculated location in the image, differ. We take the second avenue and show that the errors are small. Intuitively, this is the case since the point where the object is located in the image has no translation error (but might still have orientation errors). Other points on the model may also have translation errors, but obviously, they would be small.

We utilize vanishing points and lines for camera calibration, an approach which is best suited for architectural scenes, due to the orthogonality and parallelism that exist in these scenes. In [Faugeras et al., 1998] the problem of reconstructing a textured three dimensional model of a static urban scene viewed by one or several cameras whose relative positions are unknown, is addressed. In [Cipolla et al., 1999] the problem of recovering 3D models from two or more uncalibrated images of architectural scenes is addressed and a 3D model is recovered. In [Criminisi et al., 2000, Liebowitz et al., 1999] methods for creating 3D graphical models of scenes from one or two images are presented. The techniques presented are for metric reconstruction and for correct representation of angles and length ratios. In particular, in [Criminisi et al., 2000], it is described how 3D affine measurements can be computed from a single perspective view of a scene given only a minimal geometric information which can be determined from the image.

Our work has three novel aspects. First, we describe a technique for automatically extracting vanishing points from a given image. We show that this problem reduces to the problem of classifying lines in the image into sets associated with the principal vanishing point.

Second, we present a new theoretical error-analysis method. This method can be used for measuring the quality of various augmented reality algorithms. We also compare the method to the Monte-Carlo estimation technique, in order to demonstrate its quality.

Third, we implemented a system and tested it on many indoor scenes, planting in them various 3D models. The models were placed on various planes in the picture – on furniture, on the floor, hanging from the ceiling and hanging on walls. In all cases, the size, orientation and position of the objects were highly accurate. A system like this has a potential for various commercial applications, for instance for furniture or household companies. The client can bring an image of a room and the system can show the client suggestions for interior design of the room using 3D models of objects from the company’s stock.

The paper continues as follows. Section 2 describes our general approach. Section 3 focuses on the extraction of the principal vanishing points. Error analysis methods are described in Section 4. We view this section as the main theoretical contribution of the paper. The results of testing our system on real indoor scenes are presented in Section 5. Section 6 concludes the paper.

## 2 The Approach

This section explains the notations and describes the methods used for camera calibration. The goal is to find the parameters of the camera in order to construct a virtual camera with the same parameters as the real camera. This virtual camera projects the virtual objects on the real image.

A camera maps the 3D world to the image plane. Given a point  $P = (X, Y, Z)$  in the world its image  $p$  in homogeneous coordinates is:  $p \cong R(P - C)$ , where  $\cong$  means equal up to a multiplicative constant,  $C$  is the camera center and the orientation of the camera is represented by the rotation matrix  $R$ . The camera center and the rotation matrix are the *extrinsic* parameters of the image.

The internal calibration matrix  $K$  converts a normalized point  $p$  to its pixel coordinates:

$$K = \begin{pmatrix} f_X & 0 & u_0 \\ 0 & f_Y & v_0 \\ 0 & 0 & 1 \end{pmatrix},$$

where  $f_X$  and  $f_Y$  are the focal lengths in the  $X$  and  $Y$  directions respectively and  $(u_0, v_0)$  is the principal point (the projection of points in direction  $(0, 0, 1)$ ). These parameters are known as the *intrinsic* calibration parameters of the camera. We assume that the aspect ratio is known and thus we can solve the problem only for the case where  $f = f_X = f_Y$ .

The projection of the point  $P$  to the pixel  $s = (u, v, 1)$  in the image is:

$$s \cong KR(P - C) = [KR | -KRC] \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}.$$

We denote by  $M = [KR | -KRC]$  the  $4 \times 3$  projection matrix. Thus when we are given  $M$  we can compute the projection of any 3D point  $P$  to the image.

Finding the parameters of the virtual camera is done in three steps: extraction of the principal vanishing points, recovery of the projection matrix given the vanishing points, and recovery of the external calibration parameters. Moreover, to position a virtual object in the image, the user selects a point in the image. The world coordinates of this point need to be calculated. Below we elaborate on these four issues.

**1. Extraction of principal vanishing points:** Given an image containing lines which are parallel in the world, these lines meet at a point in the image termed a vanishing point. Three sets of lines parallel to the three principal axes are automatically extracted from the image, and their vanishing points  $V_x = (u_x, v_x, 1)$ ,  $V_y = (u_y, v_y, 1)$  and  $V_z = (u_z, v_z, 1)$  are computed. We discuss this issue in the next section.

**2. Recovery of the projection matrix:** We assume that  $o$  is the image of the origin of the world coordinate system  $O$ . This point is selected arbitrarily by the user (e.g. a corner in the room) and the world coordinate system is aligned with the principal axes in the room. The vanishing points  $V_x, V_y, V_z$  jointly with  $o$  are the images of the world points  $(1, 0, 0, 0)$ ,  $(0, 1, 0, 0)$ ,  $(0, 0, 1, 0)$  and  $(0, 0, 0, 1)$  respectively. We denote the columns of  $M$  by  $m_x, m_y, m_z, m_o$ . By applying  $M$  to these four special points we get

$$\lambda_1 V_x = m_x, \quad \lambda_2 V_y = m_y, \quad \lambda_3 V_z = m_z, \quad \lambda_4 o = m_o,$$

where the  $\lambda_i$ 's are unknown multiplicative constants.

In order to recover the unknown  $\lambda_i$ 's, recall that  $[m_x | m_y | m_z] = KR$ . Thus  $(KR)(KR)^T = KR R^T K^T = KK^T$ . Representing  $KR$  by the vanishing points we get

$$KR = \begin{pmatrix} \lambda_1 u_x & \lambda_2 u_y & \lambda_3 u_z \\ \lambda_1 v_x & \lambda_2 v_y & \lambda_3 v_z \\ \lambda_1 & \lambda_2 & \lambda_3 \end{pmatrix} = \begin{pmatrix} u_x & u_y & u_z \\ v_x & v_y & v_z \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{pmatrix}.$$

Thus

$$\begin{aligned} (KR)(KR)^T &= \begin{pmatrix} u_x & u_y & u_z \\ v_x & v_y & v_z \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \lambda_1^2 & 0 & 0 \\ 0 & \lambda_2^2 & 0 \\ 0 & 0 & \lambda_3^2 \end{pmatrix} \begin{pmatrix} u_x & v_x & 1 \\ u_y & v_y & 1 \\ u_z & v_z & 1 \end{pmatrix} \\ &= \begin{pmatrix} f^2 + u_0^2 & u_0 v_0 & u_0 \\ u_0 v_0 & f^2 + v_0^2 & v_0 \\ u_0 & v_0 & 1 \end{pmatrix} = KK^T. \end{aligned}$$

This yields six equations in the six unknowns  $f, u_0, v_0, \lambda_1, \lambda_2, \lambda_3$  which can be readily solved yielding a unique solution.

In order to complete the recovery of the projection matrix we still have to recover the last column  $m_o$  which is equal to  $\lambda_4 o$  for an unknown value of  $\lambda_4$ .  $\lambda_4$  represents the depth of the world origin from the camera center in the direction of the principal ray. In order to obtain this value we require input from the user as will be explained in Section 5.

**3. Recovery of external calibration parameters:** The rotation matrix  $R$  can be recovered by multiplying  $K^{-1}$  by  $KR$  which are both known. The location of the camera  $C$  is orthogonal to the three rows of  $M$ . Using this constraint  $C$  can be computed analytically.

**4. Calculating the world coordinates:** Since the projection matrix is not invertible, it is impossible to recover the world coordinates of a point from the image coordinates without extra information. However, if we know that the world point lies on a known plane, there is a one-to-one mapping (a homography) between the point in the image and its world coordinates. For example, for the plane  $Z = z_0$

$$s \cong M \begin{pmatrix} X \\ Y \\ z_o \\ 1 \end{pmatrix} = [m_x | m_y | m_z | m_o] \begin{pmatrix} X \\ Y \\ z_o \\ 1 \end{pmatrix} = [m_x | m_y | z_0 m_z + m_o] \begin{pmatrix} X \\ Y \\ 1 \end{pmatrix}$$

which enables us to convert image points known to lie on the plane  $Z = z_0$  to their 3D points.

### 3 Extraction of Vanishing Points

The only information that has to be extracted from the image is the set of vanishing points. This is done in three steps.

First, straight lines are extracted from the image. An edge detector (e.g., the Canny edge detector [Canny, 1986]) is first applied to the image to find its edges. Then, straight lines are

recovered from the edge map by fitting straight lines to sequences of adjacent edge elements, whose least squares matching error is less than a preset threshold.

Second, the lines are classified into four sets, one for each of the three principal axis directions and the fourth for straight lines not parallel to any of the principal axes. This process is described in Section 3.2.

Finally, for a given set of lines associated with a specific principal axis, the vanishing point is estimated. Section 3.1 describes this calculation which is also used in the line classification algorithm.

### 3.1 Numerical Methods

Given a set of lines in the image which are projections of parallel lines in three dimensions, we now describe how to find its most probable vanishing point. In an ideal setting there is only one intersection point that all the lines in the set are incident on. However, inaccurate line detection leads to several intersection points in one set of projected parallel lines. We need to find a point that is the closest to the accurate vanishing point.

Let us consider a set of projected parallel lines  $L = \{l_1, \dots, l_n\}$  where  $l_i = (a_i, b_i, c_i)$ . Recall that a point  $p = (u, v, 1)$  is incident on  $l = (a, b, c)$  if  $au + bv + c = 0$  or  $l \cdot p = 0$ . In order to find the point closest to all the lines in a least-squares sense, it is customary to solve the least squares problem of  $[A|B](u, v) = -C$  where  $A, B$  and  $C$  are the vectors of the  $a_i$ 's,  $b_i$ 's, and  $c_i$ 's respectively. This solution however, is not the *approximate maximum likelihood (AML)* estimate of the vanishing point because the above setting is equivalent to assuming that there is an error in the vanishing point and not in the lines  $l_i$ . The truth however is the opposite.

We want to find the point  $p$  that is the AML estimate of the vanishing point to achieve more accuracy [Adam et al., 2001]. We call the noise-free set of lines  $L$  and the measured set of lines  $L'$ . We say that a set of lines supports a point  $p$  if all the lines in the set meet at the point  $p$ . We want to find a set  $L''$  that supports a point  $p$  that is the most likely among all possible sets of lines, given the set of measured lines  $L'$ . In other words, we want to find the point  $p$  that maximizes the likelihood of  $L''$  being the true set of lines, given the measured set  $L'$ . Thus given a candidate point  $p$  we have to give an estimate for the probability that it is the vanishing point.

Figure 2 illustrates the geometric solution for finding the set  $L''$ . We assume that  $q'_1 = (x'_1, y'_1)$  and  $q'_2 = (x'_2, y'_2)$  are the endpoints on a line  $l'$  in  $L'$ . We want to find a line  $l''$  which is incident on  $p$  such that the squared distance from  $q'_1$  and  $q'_2$  to that line (distances to  $q''_1$  and  $q''_2$ ) is minimal. In other words, let  $d_1$  and  $d_2$  be the distances between  $q'_1$  and  $q''_1$ , and between

$q'_2$  and  $q''_2$ . Our goal is to minimize the sum of distances ( $d_1^2 + d_2^2$ ) from the new endpoints to the old endpoints. The angle  $\theta$  between the  $x$  axis and the line created by the optimal  $q'_1$  and  $q''_2$  and the distances are:

$$\begin{aligned}
\theta &= \frac{1}{2} \operatorname{atan2}(2(x'_1 - u)(y'_1 - v) + (x'_2 - u)(y'_2 - v), \\
&\quad (x'_1 - u)^2 + (x'_2 - u)^2 - (y'_1 - v)^2 - (y'_2 - v)^2) \\
d_1^2 &= (-\sin \theta (x'_1 - u) + \cos \theta (y'_1 - v))^2, \\
d_2^2 &= (-\sin \theta (x'_2 - u) + \cos \theta (y'_2 - v))^2 \\
\Delta &= d_1^2 + d_2^2.
\end{aligned} \tag{1}$$

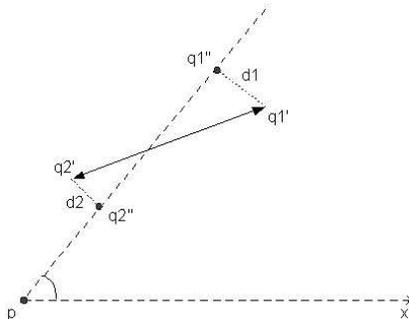


Figure 2: Illustration of the likelihood function.

The line segments of  $L'$  are independent on each other, therefore to maximize the likelihood of the set  $L''$  being the true set, the sum of distances of all lines in  $L' = (X'_1, Y'_1, X'_2, Y'_2)$  needs to be minimized. The function  $F$  to be minimized is:

$$\begin{aligned}
F(u, v, X'_1, Y'_1, X'_2, Y'_2) &= \sum_{i=1}^n \Delta_i = \\
&\sum_{i=1}^n \{(-\sin \theta_i (x'_{i1} - u) + \cos \theta_i (y'_{i1} - v))^2 \\
&\quad + (-\sin \theta_i (x'_{i2} - u) + \cos \theta_i (y'_{i2} - v))^2\}.
\end{aligned} \tag{2}$$

Finding the point  $p = (u, v)$  which minimizes this function gives the best estimate for the vanishing point we seek. To find  $p$  we apply a non-linear optimization method.

### 3.2 Line Classification

Given  $n$  lines we hereby describe how to classify them into four sets, one for each of the three principal axis directions and the fourth for straight lines not parallel to any of the principal axes. Obviously, checking all the  $4^n$  possibilities is infeasible. We therefore use a variant of

the Random Sample Consensus (RANSAC) paradigm [Fischler and Bolles, 1981] to classify the lines. In this paradigm, a minimal number of features required to estimate the model are chosen at random, the model is computed and is verified using the remaining features. This process is repeated until a correct model is found.

In our case, the model is the three vanishing points and the features are three pairs of lines. Thus, three pairs of lines are chosen at random and it is hypothesized that each pair of lines is incident on a different vanishing point. Under this assumption the three intersection points are computed and are hypothesized to be the principal vanishing points (obviously due to measurement errors these lines do not intersect at the true vanishing points). The remaining lines are assumed to pass at a random distance from them. If the hypothesis were correct, many lines would have passed close to one of the principal vanishing points.

Given an hypothesis  $H$ , our goal is to compute a score  $s(H)$  for it such that the lower the score the higher the probability that the hypothesis is correct. The score  $s(H)$  is defined as the sum of scores computed for all the lines recovered from the image.

For each line  $l'_i$ ,  $1 \leq i \leq n$ , the values  $\Delta_{i1}, \Delta_{i2}, \Delta_{i3}$ , as defined in Equation 1, are computed for the three candidate vanishing points. Each value  $\Delta_{ij}$  is proportional to the  $-\log$  of the probability that line  $l'_i$  is incident on vanishing point  $j$  for  $1 \leq j \leq 3$ . Line  $l'_i$  is associated with the vanishing point having the smallest of the three values  $\Delta_{min} = \min(\Delta_{i1}, \Delta_{i2}, \Delta_{i3})$ . If  $\Delta_{min}$  is larger than a preset threshold  $D$  we assign the value  $D$  to this line which means that the line is probably not incident on any of the three vanishing points and therefore we are indifferent to its actual distances to the vanishing points.

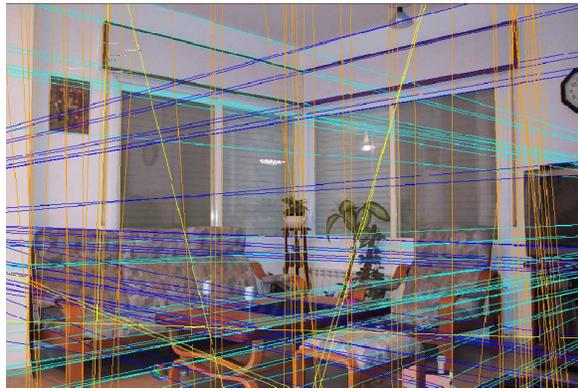
This random selection is repeated many times (in practice, several thousand times) and the hypothesis with the smallest score is chosen. This score represents the quality of the hypothesis in two ways. First, the larger the number of lines which are associated with a vanishing point (few lines got the score  $D$ ), the better the score is. This feature causes hypotheses for which each pair of lines belongs to a different vanishing point to get a low score. Second, among all correct hypotheses (i.e., each pair of lines belongs to a different principal vanishing point), the hypothesis chosen more accurately estimates the three vanishing points. This results from the properties of the AML estimate we are using. This method (Equations 1 and 2) achieves the minimal sum of squared distances by definition, thus, hypotheses having lower scores yield more accurate estimates.

Figure 3 illustrates an example of an image on which this procedure was run. Figure 3(a) shows the results of the straight line detection process. Figure 3(b) shows the lines that have been classified as belonging to each of the three vanishing points (in blue, cyan and orange) and lines which do not belong to any of them (in yellow). In this figure the vanishing points

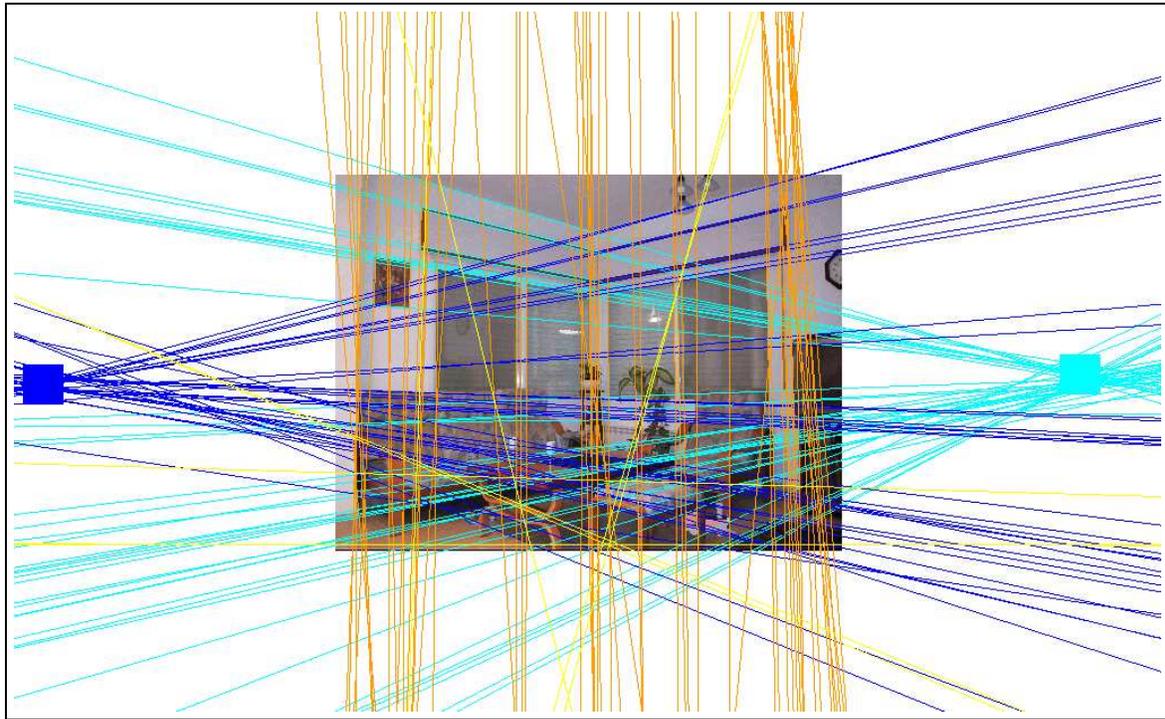
are not visible. To see them, Figure 3(c) presents a zoomed out version of the image in which two of the vanishing points are visible.



(a) Results of the edge/line detector



(b) The classification results



(c) A zoomed out image of (b) in which two of the three vanishing points are visible

Figure 3: Results of the line classification algorithm: The blue, cyan and orange lines belong to the three principal vanishing points and the yellow lines are unclassified lines.

## 4 Error Estimation and Analysis

This section focuses on the accuracy of the projection of 3D virtual objects onto a single 2D image. The quality of the recovered camera calibration determines the accuracy of the

parameters of the virtual camera, used to “take a picture” of the virtual objects. Two methods for estimating the errors are described: a novel analytical method and a Monte Carlo method. Both methods compute the covariance matrices of the parameters, which give us the distribution of the errors of the computed values as functions of the errors in the measured values, thus estimating the real errors. We will show that our proposed analytical method is almost as accurate as the Monte Carlo method. Since it is much faster, it can be the method of choice for augmented reality algorithms.

## 4.1 Analytical Estimation of Covariance

We describe below an analytical method for computing the covariance matrices of several parameters. Our objective is to estimate the dispersion of the vanishing points, the principal point, the focal length, the camera location in the world coordinate system and the projection of points in the world to the image.

Let the ideal input be an  $N \times 1$  vector  $X$ , the observed perturbed input be an  $N \times 1$  vector  $\hat{X} = X + \Delta X$ , where  $\Delta X$  is a random vector of noise. Let the correct parameters sought be a  $K \times 1$  vector  $\Theta$ , and the calculated parameters be a randomly perturbed  $K \times 1$  vector  $\hat{\Theta} = \Theta + \Delta\Theta$ , where  $\Delta\Theta$  is a random perturbation on  $\Theta$  induced by the random perturbation  $\Delta X$  on  $X$ .

We base our estimate on a first order approximation derived in [Haralick, 1996], describing how to propagate approximately additive random perturbations through an algorithm step. Covariance propagation allows us to determine the covariance of  $\Theta$  where  $\Theta = G(X)$  from the covariance of  $X$ . We assume that  $G$  is a non-linear function that is approximately linear in the vicinity of the measured values of  $X$ . In this case the covariance matrix of  $\Theta$  is given by  $J\Sigma J^T$ , where  $J$  is the Jacobian matrix of  $G$  evaluated at  $X$ .

A related problem is when  $G$  is not given explicitly. Instead, a function  $F$ ,

$$\min_{\Theta} F(\Theta, X)$$

relates the vectors  $X$  and  $\Theta$  and the vectors  $X + \Delta X$  and  $\Theta + \Delta\Theta$ .  $F$  has finite first and second partial derivatives with respect to each component of  $X$  and  $\Theta$  including all the second partial derivatives. The basic problem is thus: given  $\hat{X} = X + \Delta X$ , determine  $\hat{\Theta} = \Theta + \Delta\Theta$  that minimizes  $F(\hat{X}, \hat{\Theta})$  assuming that  $\Theta$  minimizes  $F(X, \Theta)$ . Let

$$g(\Theta, X) = \frac{\partial F}{\partial \Theta}(\Theta, X).$$

Since  $\hat{\Theta}$  minimizes  $F(\hat{X}, \hat{\Theta})$  and  $\Theta$  minimizes  $F(X, \Theta)$ , both  $g(\hat{X}, \hat{\Theta})$  and  $g(X, \Theta)$  equal

zero. By taking a Taylor series expansion of  $g$  around  $(X, \Theta)$  we obtain a first order approximation:

$$\Delta\Theta = - \left( \frac{\partial g}{\partial \Theta}(\Theta, X) \right)^{-1} \frac{\partial g}{\partial X}(\Theta, X) \Delta X \quad (3)$$

This relation states how the random perturbation  $\Delta X$  on  $X$  propagates to the random perturbation  $\Delta\Theta$  on  $\Theta$ . If  $E[\Delta X] = 0$  then from Equation 3,  $E[\Delta\Theta] = 0$  to a first order approximation. The covariance matrix of  $\hat{\Theta}$  is then:

$$\Sigma_{\hat{\Theta}} = \mathcal{M} \Sigma_{\hat{X}} \mathcal{M}^T, \quad (4)$$

where

$$\Sigma_{\hat{X}} = E[\Delta X \Delta X^T], \quad \mathcal{M} = \left( \frac{\partial g}{\partial \Theta} \right)^{-1} \frac{\partial g}{\partial X}. \quad (5)$$

We can apply the above theory to our calibration calculations. Recall that our inputs are three sets of end points of image straight lines, in the three principal directions. Let  $X$  be the set of accurate end points of the lines in one direction,  $X = (x_1, y_1, x_2, y_2, \dots, x_p, y_p)^T$ , and let  $\hat{X}$  be the corresponding noisy set measured in the image,  $\hat{X} = (\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \dots, \hat{x}_p, \hat{y}_p)^T$ . The calibration algorithm first finds the principal vanishing points. In this case the function  $F$  is the function from Equation 2,  $\Theta = (u, v)$  is the accurate vanishing point and  $\hat{\Theta} = (\hat{u}, \hat{v})$  is the estimated vanishing point calculated from the noisy inputs. We can calculate the matrix  $\mathcal{M}$  in Equation 5 from:

$$\begin{aligned} g &= \begin{pmatrix} Fu \\ Fv \end{pmatrix} \\ \frac{\partial g}{\partial \Theta} &= \begin{pmatrix} F_{uu} & F_{uv} \\ F_{vu} & F_{vv} \end{pmatrix} \\ \frac{\partial g}{\partial X} &= \begin{pmatrix} F_{u\hat{x}1} & F_{u\hat{y}1} & \cdots & F_{u\hat{y}p} \\ F_{v\hat{x}1} & F_{v\hat{y}1} & \cdots & F_{v\hat{y}p} \end{pmatrix}, \end{aligned} \quad (6)$$

and use it to compute  $\Sigma_{\hat{\Theta}_{vp}}$  using Equation 4.

The next steps in the calibration process are calculations of the various parameters using deterministic non-linear functions of the vanishing points. These are the intrinsic parameters of the camera (i.e., the principal point and the focal length). In these cases the covariance matrices are obtained by:

$$\Sigma_{\hat{P}} = J \Sigma_{\hat{\Theta}} J^T, \quad (7)$$

where  $\hat{P}$  is the estimation of the required parameter and  $J$  is the Jacobian matrix of the function evaluated at the correct vanishing points.

We can use the covariance matrix estimated for the vanishing points to estimate the covariance matrix of the principal point and the covariance matrices of the scale factors  $\lambda_1, \lambda_2, \lambda_3$  and  $f$ .

The next step in the calibration phase is to find the location of the camera in the world coordinate system. The elements of the projection matrix are functions of the principal vanishing points and of the scale factors. They also depend on a couple of parameters given by the user: the global scale factor  $\lambda_4$  and the image coordinates of the world origin  $(o_x, o_y)^T$ . Using Equation 7 the covariance of the estimated location of the camera can be calculated as well.

After computing the camera's intrinsic and extrinsic parameters we can calculate the projection of points in 3D. The system computes the world coordinates of a point marked by the point in the image where a virtual object should be placed. Estimating the error of the projection of the virtual object is more difficult than the previous estimations because it consists of two phases. We first estimate the covariance matrices of the components of the projection matrix as functions of the principal vanishing points, and then use the recovered covariance matrices to estimate the covariance matrices of the projected points to the image. This is done as follows:

$$u = \frac{M^{1T}(X, Y, Z, 1)^T}{M^{3T}(X, Y, Z, 1)^T}, \quad v = \frac{M^{2T}(X, Y, Z, 1)^T}{M^{3T}(X, Y, Z, 1)^T}, \quad (8)$$

where  $M^{1T}, M^{2T}$  and  $M^{3T}$  are the first, second and third rows of the projection matrix  $M$  respectively,  $(X, Y, Z, 1)^T$  is the world coordinates of a vertex of the virtual object and  $(u, v)^T$  is its image coordinates. Let  $M$  be the accurate projection matrix and  $\hat{M}$  be the estimated projection matrix. According to Equation 7 we can calculate the covariance of the estimated projection matrix as:

$$\Sigma_{\hat{M}} = J_M \Sigma_{\hat{\theta}_{vp}} J_M^T, \quad (9)$$

where  $J_M$  is the Jacobian matrix of the projection matrix evaluated at the correct vanishing points and  $\Sigma_{\hat{\theta}_{vp}}$  is the covariance matrix of the vanishing points. We can now calculate the covariance of the estimated image coordinate as:

$$\Sigma_{(\hat{u}, \hat{v})} = J_{uv} \Sigma_{\hat{M}} J_{uv}^T, \quad (10)$$

where  $J_{uv}$  is the Jacobian matrix of the image coordinates evaluated at the correct projection matrix.

## 4.2 Monte Carlo Estimation of Covariance

The Monte Carlo method estimates the covariance by exhaustive simulation [Hartley and Zisserman, 2000]. It can be used to validate the results of the analytical method. This method is also preferable in cases where the first order approximation of the covariance is not accurate enough and thus the analytical estimation may be incorrect. The goal is to statistically compute the covariance matrices of the estimated parameters, having noisy input.

We created a synthetic 3D scene which was rendered from five different view points creating five different images, as shown in Figure 4. The location of the camera between each pair of the first three images is translated while the orientation of the camera remains the same. In the last three images the orientation of the camera changes while the location of the camera remains the same.

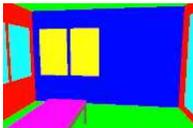
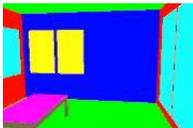
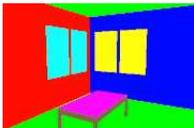
	Test1	Test2	Test3	Test4	Test5
<b>Image</b>					
<b>Rotation</b>	$(-6.5, 0, 8.0)$	$(-6.5, 0, 8.0)$	$(-6.5, 0, 8.0)$	$(-6.5, 0, 15.0)$	$(-6.5, 0, 30.0)$
<b>Center</b>	$(-200, -410, 105)$	$(-250, -430, 105)$	$(-300, -460, 105)$	$(-300, -460, 105)$	$(-300, -460, 105)$

Figure 4: The synthetic test images, the angles ( $\circ$ ) between viewing direction and the principal directions and the camera's center in the world (cm).

The inputs to each test are the following: 1) The intrinsic parameters of the camera. 2) Three sets of end points of eight straight lines, each belonging to one set of parallel lines in the world coordinate system: X, Y and Z. The end points are given as 3D points and are projected to the image by the computed projection matrix. 3) The location of the world origin in the image. 4) One additional point in the image and its 3D distance from the origin of the world coordinate system. This information is used for calculating the global scale factor. 5) Eight 3D points that are used as test points for checking the accuracy of the virtual object registration. These points are translated to 2D image points by the correct projection matrix once at the beginning of the test. Each point is accompanied by information regarding the plane the point lies on and its distance from the principal plane which is parallel to it. See Figure 5.

First, the correct parameters are computed using the pure inputs. Afterward, random noise is added to each end point of the input straight lines. The noise has a Gaussian distribution with zero mean and standard deviation values  $\sigma = 0.2, 0.4, 0.6, 0.8, 1$  (in pixels).

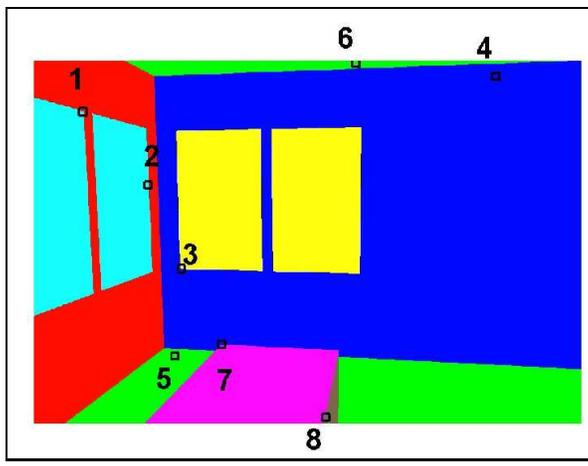


Figure 5: Test1 image with the 8 test points marked on it.

We conducted 1000 experiments with each level of  $\sigma$ .

For each set of noisy inputs we compute the principal vanishing points, the intrinsic parameters of the camera, the projection matrix and the rotation and the translation of the camera. We also calculate the virtual object registration in the image, by computing the projection of 3D points onto the 2D image, using the computed projection matrix. In our setup, the user marks a point in the image where the virtual object should be placed. This 2D point is first transformed to a 3D point using the computed projection matrix, and the virtual object is placed in that coordinate in 3D. The object's vertices are then projected to the image according to the computed projection matrix.

In our experiments we created a test virtual object with 18 vertices, where the distances of these vertices from the center of the virtual object are 30cm, 60cm and 90cm in the directions  $X$ ,  $Y$  and  $Z$  in the world coordinate system. These 18 test vectors are projected to the image using the computed projection matrix.

We compare the resulting parameters to the pure parameters and calculate the errors as follows: 1) The distance between the correct vanishing points to the perturbed vanishing points in pixels. 2) The distance from the pure principal point to the noisy principal point in pixels in each axis, divided by the size of the image in each axis. 3) The relative error in the estimated focal length. 4) In order to calculate the error in the rotation of the camera, the computed rotation matrix is multiplied by the inverse of the correct rotation matrix. The result is the rotation between the pure orientation of the camera and the perturbed orientation of the camera. We compute the quaternion of this rotation matrix to get its rotation angle and divide it by the field of view to get the relative error. 5) The error in the position of the camera center in centimeters. 6) The error in the image of the projected 3D

points, represented by the distance and the angle between the correct projected vector and its corresponding noisy vector.

After performing 1000 simulations at a specific noise level we calculate the mean and the standard deviation of each kind of error. The standard deviation is estimated as the  $2/3 * 1000 = 666$  element in the sorted vector of a specific error (which is a robust estimate for the standard deviation).

### 4.3 Error Estimation Results

In this section we present the results of the analytical error estimation as well as the results of the Monte Carlo error estimation, and compare the two for the test images in Figure 4.

Table 1 shows the average of the errors of the intrinsic and the extrinsic parameters of the camera for all the test images calculated by the Monte Carlo error estimation method for noise level of  $\sigma = 1$ . Table 2 shows the standard deviation of the errors of the intrinsic and the extrinsic parameters of the camera for the test images calculated by the Monte Carlo error estimation method. It can be seen that the errors are negligible relative to their standard deviations. Thus our assumption of zero mean, made in the analytical method, is justified. The results for test4 and test5 are better than the others, due to the fact that in these images the viewing direction of the camera has a large angle with two of the principal directions. We conclude that the best camera orientation is when the world three planes  $XY$ ,  $XZ$  and  $YZ$  form an angle close to  $\pi/4$  with the image plane. In this case the vanishing points are closer to the image and the accuracy of their localization in the image increases [Caprile and Torre, 1990].

Error in	Test1	Test2	Test3	Test4	Test5
Focal length % of focal length	0.066	0.051	0.088	0.033	0.002
u0 % of X dim of image	0.016	0.023	-0.020	-0.002	0.034
v0 % of Y dim of image	-0.007	-0.016	-0.006	-0.006	-0.011
Camera center X coord. in cm	-0.596	-0.836	0.451	-1.502	-0.29
Camera center Y coord. in cm	-0.233	0.085	1.281	-0.849	0.55
Camera center Z coord. in cm	0.465	0.391	0.619	0.447	0.547

Table 1: The average errors using the Monte Carlo error estimation method for noise level of  $\sigma = 1$ .

Figure 6 shows the results of the estimation of the errors of the intrinsic parameters of the camera and the translation of the camera using both methods. The graphs show that both methods give similar results, justifying the assumption made in the analytical method, that

Error in	Test1	Test2	Test3	Test4	Test5
Focal length % of focal length	4.4306	4.1744	4.6565	2.4134	1.0077
u0 % of X dim of image	1.0195	1.0161	1.1205	1.0926	1.1864
v0 % of Y dim of image	1.365	1.2159	1.3124	0.721	0.5257
Rotation angle	0.4175	0.3937	0.4277	0.3839	0.3978
Camera center	19.084	18.801	22.521	12.498	6.5957

Table 2: The standard deviation of the errors for noise level  $\sigma = 1$ .

the behavior is linear.

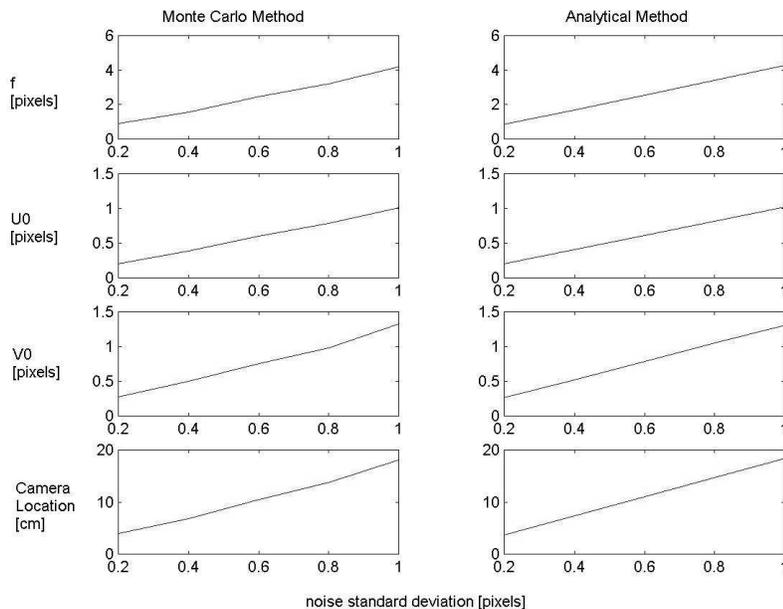


Figure 6: The standard deviation of the errors in test1 as a function of the input noise level. Left: Monte Carlo. Right: analytical.

As described above, 8 points in the synthetic scene were chosen (Figure 5) and the projection of 18 vectors originating at each test point were calculated. The angle between each of the pure 18 vectors and the perturbed 18 vectors, is calculated. Table 3 summarizes the results. It shows that the angles are very small,  $0.072^\circ$  on average, thus the perturbed vectors and the pure vectors are approximately in the same direction in the image.

We also calculated the distance in pixels between each of the pure 18 points (for each test point) and the perturbed 18 points, divided by the length of the pure projection of each vector. We conclude that the distortion of the virtual objects projected to the image is un-noticeable to the viewer, being at most one or two pixels for standard household items.

		Test1	Test2	Test3	Test4	Test5
test point 1 (0,-110,50)	dist. from cam. [cm]	364.73	409.79	464.25	464.25	464.25
	error in length %	0.282	0.306	0.358	0.368	0.480
	error in angle [deg]	0.027	0.021	0.018	0.022	0.068
test point 2 (0,-20,100)	dist. from cam. [cm]	438.32	480.23	532.56	532.56	532.56
	error in length %	0.275	0.290	0.329	0.356	0.508
	error in angle [deg]	0.035	0.029	0.025	0.031	0.098
test point 3 (-20,0,180)	dist. from cam. [cm]	454.01	493.38	543.71	543.71	543.71
	error in length %	1.652	1.565	1.739	1.079	0.861
	error in angle [deg]	0.063	0.063	0.060	0.069	0.153
test point 4 (-300,0,10)	dist. from cam. [cm]	432.58	443.20	469.71	469.71	469.71
	error in length %	3.877	3.680	4.149	2.264	1.296
	error in angle [deg]	0.149	0.110	0.099	0.113	0.246
test point 5 (-20,-20,260)	dist. from cam. [cm]	456.65	495	544.08	544.08	544.08
	error in length %	5.474	5.033	5.518	2.928	1.586
	error in angle [deg]	0.092	0.072	0.070	0.063	0.133
test point 6 (-185,-20,0)	dist. from cam. [cm]	404.17	428.19	466.74	466.74	466.74
	error in length %	2.299	2.091	2.227	1.261	0.774
	error in angle [deg]	0.057	0.049	0.051	0.050	0.087
test point 7 (-100,- 120,210)	dist. from cam. [cm]	324.23	360.03	408.20	408.20	408.20
	error in length %	0.249	0.232	0.248	0.233	0.270
	error in angle [deg]	0.028	0.028	0.029	0.028	0.042
test point 8 (-180,- 210,210)	dist. from cam. [cm]	226.77	253.62	296.52	296.52	296.52
	error in length %	0.255	0.245	0.259	0.242	0.260
	error in angle [deg]	0.123	0.067	0.040	0.056	0.224

Table 3: The standard deviation of the errors of the projected test vectors for images with noise with  $\sigma = 1$ (pixel). The first row is the distance of the point from the camera. The second row is the average standard deviation of the errors in the lengths of the projected vectors, relative to their lengths. The third row is the average standard deviation of the errors in the angles between the projected and the pure vectors.

In conclusion, the error results obtained by both methods are similar except for the estimation of the covariance matrices of the projection of the 3D points, where the differences between the estimates is about 50% – 65%. This is due to the first order approximation made by the analytical method. These estimates have proven that a system built using our approach will yield sufficiently accurate results.

## 5 Implementation and Results

We built an interactive system that lets users augment a single indoor image with virtual objects. The system is implemented on a PC with the Windows NT operating system, using OpenGL. The test images were taken by a Canon Powershot S100 digital camera.

The line classification algorithm described in Section 3 is first applied to the image and creates four sets of lines:  $X$  lines,  $Y$  lines,  $Z$  lines and non-aligned lines. From the first three sets the vanishing points are estimated.

In order to calculate the extrinsic parameters of the camera, the user needs to supply the system with two inputs: the location of the origin of the world in the image and a real-world distance (in cm) between two points in the image. For example, in Figure 7 the cyan dot is the world origin and the length of the cyan line in the world is given by the user.



Figure 7: The user interactively specifies the cyan dot and the cyan line and its length in the world.

To insert a virtual object into the image, the user needs to specify its desired location on the image with a mouse click, and then set the distance of the plane on which the object should be placed from the principal plane parallel to it. After placing the object, the user can rotate it around its local rotation axis. Note that for a plane in a general position, the

user has to specify its 3D parameters.

Next, as described in Section 2, the system computes five transformation matrices and applies them to the object as follows:

$$Projected\_Object \cong [K][R][T][GR][LR](Virtual\_Object).$$

[K] is the camera's intrinsic parameters matrix. [R] is the camera's rotation matrix. [T] is the translation of the virtual object to its location in the world coordinate system relative to the location of the camera. [GR] is the global rotation matrix; it is computed according to the surface which the virtual object touches. Finally, [LR] is the local rotation matrix which rotates the virtual object in its local coordinate system.

We implemented a semi-automatic capability in the system that deals with occlusions of virtual objects by real objects. *Intelligent scissors* are used to let the user define *blobs* in the image. Each blob represents a closed area in the image. Specifically, we define three types of blobs. A *plane blobs* includes pixels which belong to a plane in the world. A *general blob* defines any closed area. Both types of blobs can occlude virtual objects. Finally, a *back blob* is a region in the image which is placed behind all objects, both virtual and real, thus cannot occlude objects (e.g., walls). The user can interactively reorder blobs by pushing them backward or popping them forward. In Figure 8, the user placed a toy train in the living room. Without dealing with occlusion the toy seems to be floating in the air as shown in Figure 8(a). By popping the foot-stool blob forward, the train is positioned correctly (Figure 8(b)).



(a)



(b)

Figure 8: (a) A toy train is placed in the scene incorrectly occluding the foot-stool. (b) The foot-stool blob is popped forward yielding a correct result.

Figure 9 reinforces the conclusions drawn in Section 4.3 regarding the accuracy of the registration. We created a rectangular virtual object and inserted it into a real image in a

couple of locations. Its height is equal to the height of the white area it is placed on. It can be seen that the virtual objects have accurate heights and correct orientations.



Figure 9: Two rectangular objects are placed accurately in the image.

Figures 11–12 illustrate some results obtained by our system. The images show three different rooms and they were taken from different viewpoints. Each image is augmented with several virtual objects, some hanging from the ceiling, others are placed on the floor or on the table, and some are hanged on the walls, hence demonstrating the use of the system on various planes.

## 6 Conclusion

This paper has addressed the problem of camera calibration recovery and virtual object registration given a single image of an indoor architectural scene. Our approach has exploited the orthogonality and the parallelism which are prevalent in such scenes. The method is based on finding the principal vanishing points and using them to construct the projection matrix, which encapsulates both the intrinsic and the extrinsic parameters of the camera. With these parameters, we can construct a virtual camera which has the same parameters as the real camera. The virtual camera lets us augment the image by overlaying virtual objects upon the real image in the correct position and orientation.

A major contribution of this paper is the error analysis of the camera calibration and the object registration method. We have described an analytical method for computing the covariance matrices of the camera’s parameters and of the virtual object registration in the real image. We have also described a Monte Carlo estimation of the same covariance matrices and compared them to the covariance matrices computed by the analytical method. The



Figure 10: An augmented family room: the wine glasses, vases, candle holders, letter holder, and the shoe are all artificial. Note that the two vases have different sizes due to perspective.



Figure 11: An augmented study: The letter holder, the phone (both hung on the wall, partially occluded by the computer), the green book, the tea cups, the wine glass, the bottle, the toy train and the vase, are all artificial.



Figure 12: An augmented living room (excuse the virtual mess): A shoe and a toy train are placed on the floor, both partially occluded. All the objects on the table are artificial as well.

analytical method is preferable in estimating the covariance matrices of the camera parameters because it is more computationally efficient. The Monte Carlo method, however, is preferable in estimating the covariance matrices of the object registration errors, because it is more accurate.

We have shown that the standard deviations of the errors of the intrinsic and the extrinsic parameters of the camera are low enough to be useful for augmenting a single image with virtual objects. We have also shown that the distortion of virtual objects is un-noticeable to the viewer. Finally, we have developed a system that implements camera calibration recovery and object registration, while requiring minimal information from the end-user. Placing virtual objects in the scene is done without constructing a full three-dimensional model of the real scene.

Future work includes searching for methods for calibrating images with less than three principal vanishing points, and integration of our system with known methods for handling occlusion, lighting and shading. In addition, the only time-consuming task is line detection and classification. Thus, to make the system run in real-time, the only change necessary is to track the lines rather than finding them independently at every frame.

## References

- [Adam et al., 2001] Adam, A., Rivlin, E., and Shimshoni, I. (2001). Computing the sensory uncertainty field of a vision-based localization sensor. *IEEE Trans. on Robotics and Automation*, 17(3):258–267.
- [Antone and Teller, 2001] Antone, M. and Teller, S. (2001). Automatic recovery of relative camera rotations in urban scenes. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, pages II-282–289.
- [Azuma, 1997] Azuma, R. T. (1997). A survey of augmented reality. *Presence: Teleoperators and Virtual Environments*, 6(4):355–385.
- [Berger et al., 1999] Berger, M.-O., Wrobel-Dautcourt, B., Petitjean, S., and Simon, G. (1999). Mixing synthetic and video images of an outdoor urban environment. *Machine Vision Applications*, 11(3):145–159.
- [Boivin and Gagalowicz, 2001] Boivin, S. and Gagalowicz, A. (2001). Image-based rendering of diffuse, specular and glossy surfaces from a single image. In *Proc. ACM SIGGRAPH*, pages 107–116.

- [Canny, 1986] Canny, J. (1986). A computational approach to edge detection. *IEEE Trans. Patt. Anal. Mach. Intell.*, 8(6):679–698.
- [Caprile and Torre, 1990] Caprile, B. and Torre, V. (1990). Using vanishing points for camera calibration. *Int. J. of Comp. Vision*, 4(2):127–140.
- [Chang et al., 1999] Chang, C., Bishop, G., and Lastra, A. (1999). Ldi tree: A hierarchical representation for image-based rendering. In *Proc. ACM SIGGRAPH*, pages 291–298.
- [Cipolla et al., 1999] Cipolla, R., Robertson, D., and Boyer, E. (1999). Photobuilder - 3D models of architectural scenes from uncalibrated images. In *IEEE International Conference on Multimedia Computing and Systems*, pages 25–31.
- [Criminisi et al., 2000] Criminisi, A., Reid, I., and Zisserman, A. (2000). Single view metrology. *Int. J. of Comp. Vision*, 40(2):123–148.
- [de Agapito et al., 1999] de Agapito, L., Hartley, R., and Hayman, E. (1999). Linear calibration of a rotating and zooming camera. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, pages I:15–21.
- [Debevec, 1998] Debevec, P. (1998). Rendering synthetic objects into real scenes: Bridging traditional and image-based graphics with global illumination and high dynamic range photography. In *Proc. ACM SIGGRAPH*, pages 189–198.
- [Faugeras et al., 1998] Faugeras, O., Laveau, S., Robert, L., Csurka, G., and Zeller, C. (1998). 3-D reconstruction of urban scenes from image sequences. *Comp. Vis. Im. Understanding*, 69(3):292–309.
- [Fischler and Bolles, 1981] Fischler, M. and Bolles, R. (1981). Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Comm. of the ACM*, 24(6):381–395.
- [Haralick, 1996] Haralick, R. (1996). Propogating covariance in computer vision. *Int. J. of Patt. Recog. and Art. Intell.*, 10:561–572.
- [Hartley and Zisserman, 2000] Hartley, R. and Zisserman, A. (2000). *Multiple Views Geoemtry in Computer Vision*. Cambridge University Press.
- [Kato and Billinghurst, 1999] Kato, H. and Billinghurst, M. (1999). Marker tracking and hmd calibration for a video-based augmented reality conferencing system. In *Second International Workshop Augmented Reality*, pages 85–94.

- [Koudelka et al., 2001] Koudelka, M., Belhumeur, P., Magda, S., and Kriegman, D. (2001). Image-based modeling and rendering of surfaces with arbitrary brdfs. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, pages I:568–575.
- [Kutulakos and Vallino, 1998] Kutulakos, K. and Vallino, J. (1998). Calibration-free augmented reality. *IEEE Transactions on Visualization and Computer Graphics*, 4(1):1–20.
- [Liebowitz et al., 1999] Liebowitz, D., Criminisi, A., and Zisserman, A. (1999). Creating architectural models from images. In *Proc. Eurographics*, pages 39–50.
- [MacIntyre and Coelho, 2002] MacIntyre, B. and Coelho, E. (2002). Estimating and adapting to registration errors in augmented reality systems. In *IEEE Virtual Reality Conference*, pages 924–28.
- [Mellor, 1995] Mellor, J. (1995). Real-time camera calibration for enhanced reality visualization. In *Proc. Computer Vision, Virtual Reality and Robotics in Medicine*, pages 471–475.
- [Poupyrev et al., 2002] Poupyrev, I., Tan, D., Billingham, Kato, H., Regenbrecht, R., and Tetsutani, N. (2002). Developing a generic augmented-reality interface. *IEEE Computer*, 35(3):44–50.
- [Ramamoorthi and Hanrahan, 2001] Ramamoorthi, R. and Hanrahan, P. (2001). A signal-processing framework for inverse rendering. In *Proc. ACM SIGGRAPH*, pages 117–128.
- [Shade et al., 1998] Shade, J., Gortler, S., He, L., and Szeliski, R. (1998). Layered depth images. In *Proc. ACM SIGGRAPH*, pages 231–242.
- [Teller et al., 2001] Teller, S., Antone, M., Bodnar, Z., Bosse, M., Coorg, S., Jethwa, M., and Master, N. (2001). Calibrated, registered images of an extended urban area. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*
- [Wloka and Anderson, 1995] Wloka, M. and Anderson, B. (1995). Resolving occlusion in augmented reality. In *Symp. on Interactive 3D Graphics*, pages 5–12.

# List of Figures

1	An augmented room . . . . .	3
2	Illustration of the likelihood function. . . . .	9
3	Results of the line classification algorithm: The blue, cyan and orange lines belong to the three principal vanishing points and the yellow lines are unclassified lines. . . . .	11
4	The synthetic test images, the angles ( $\phi$ ) between viewing direction and the principal directions and the camera's center in the world (cm). . . . .	15
5	Test1 image with the 8 test points marked on it. . . . .	16
6	The standard deviation of the errors in test1 as a function of the input noise level. Left: Monte Carlo. Right: analytical. . . . .	18
7	The user interactively specifies the cyan dot and the cyan line and its length in the world. . . . .	20
8	(a) A toy train is placed in the scene incorrectly occluding the foot-stool. (b) The foot-stool blob is popped forward yielding a correct result. . . . .	21
9	Two rectangular objects are placed accurately in the image. . . . .	22
10	An augmented family room: the wine glasses, vases, candle holders, letter holder, and the shoe are all artificial. Note that the two vases have different sizes due to perspective. . . . .	23
11	An augmented study: The letter holder, the phone (both hung on the wall, partially occluded by the computer), the green book, the tea cups, the wine glass, the bottle, the toy train and the vase, are all artificial. . . . .	24
12	An augmented living room (excuse the virtual mess): A shoe and a toy train are placed on the floor, both partially occluded. All the objects on the table are artificial as well. . . . .	25