

# Towards 3D Object Model Acquisition and Recognition using 3D Affine Invariants

Sven Vinther and Roberto Cipolla

Department of Engineering  
University of Cambridge  
Cambridge, CB2 1PZ, England

## Abstract

We evaluate the power of 3D affine invariants in an object recognition scheme. These invariants are actively calculated by the real-time tracking of 2D image features (*corners*) over an image sequence. This is done optimally by using a Kalman filter. Object information is located in a hash table where it is stored and retrieved using the invariants as stable indices. Recognition takes place when significant evidence for a particular shape has been found from the table. Preliminary results with real data are presented, and some of the noise problems arising due to the weak perspective approximation and corner localisation errors are discussed.

## 1 Introduction

Three dimensional object recognition is a rapidly expanding field within computer vision. Central to object recognition systems are questions of model acquisition and representation, feature extraction and matching. All of these issues are closely related, as the type of feature extraction performed will depend directly on the model adopted, and will also affect the matching technique selected. Extensive surveys describing methods of object representation and matching have been made by Chin and Dyer [1], Besl and Jain [2].

A number of working recognition schemes have been proposed for simple 3D objects. Lowe [3] presents a 3D object recognition system which uses a single image, *perceptual groupings* and viewpoint consistency constraints to detect 3D objects from 2D data. Thompson and Mundy [4] use *vertex-pairs* to derive the affine transformation between a 3D polyhedral object model and its projection into the image viewplane.

A more common approach is the use of 3D data extracted from range data or multiple views. The approach by Grimson and Perez [5] operates by examining all hypotheses between segmented range data and model surfaces, and efficiently discarding inconsistent ones by using local constraints. This is achieved by a depth first search of their *interpretation tree*, which at each level attempts to match an image feature to all model features. The pruning of subtrees after a failed match and the use of a search cutoff keeps search time to acceptable levels.

An alternative approach is *geometric invariants* [6]. Several successful 2D object recognition systems have been implemented using invariants [7]. In this paper we evaluate the power of 3D affine invariants in an object recognition scheme.

## 2 Theoretical Framework

The image of a 3D object will appear different depending on viewpoint because its 2D projection suffers from distortions, self-occlusion and a loss of depth. Solid object descriptors such as edge length, surface area and angles between lines which are invariant in 3D space, are no longer so in the 2D image. However a closer examination of the effects of projective transformations reveals that it is possible to find properties of combinations of features which remain *invariant* in different views. For example the *cross-ratio* of four points on a line remain invariant to general projection, and simpler invariants can be found for simpler transformations. A more detailed overview of this field of research and its applications can be found in [6].

### 2.1 Weak Perspective

The full perspective transformation of world points to image points  $(u, v)$  can be represented by (1). It combines the effects of aligning a camera-centered frame with a general world frame, and then projecting the points using the *pinhole camera* model.

$$\begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} t_{11} & t_{12} & t_{13} & t_{14} \\ t_{21} & t_{22} & t_{23} & t_{24} \\ t_{31} & t_{32} & t_{33} & t_{34} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \quad (1)$$

A simplification of this transformation is possible under *weak* perspective [8] (orthographic projection plus a scale factor). Weak perspective is a good approximation of full perspective when the distance from object to camera,  $Z$ , is much greater than the extent of the object  $\Delta Z$ . Under these conditions scaling factor  $s = t_{31}X_w + t_{32}Y_w + t_{33}Z_w + t_{34}$ , becomes approximately constant  $s = t_{34}$  in effect linearising the transformation.

$$\begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} t_{11} & t_{12} & t_{13} & t_{14} \\ t_{21} & t_{22} & t_{23} & t_{24} \\ 0 & 0 & 0 & t_{34} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \quad (2)$$

There are just eight degrees of freedom for this weak perspective transformation so only four points are required to evaluate the transformation, and no camera calibration is required.

### 2.2 Affine Invariants from Multiple Views

The linearity of weak perspective allows us to recover invariant 3D object structure from two or more images [6, 9, 10]. Given five 3D points  $\mathbf{X}_i, i \in \{0, \dots, 4\}$  it is possible to calculate 3 invariants [6]. This can be shown by considering a basis  $\mathbf{E}_i = \mathbf{X}_i - \mathbf{X}_0, i \in \{1, 2, 3\}$  in 3D space. Any 5th point in this basis can be expressed as a linear combination of  $\mathbf{E}_i$ 's (3) (see Figure 1).

$$\mathbf{X}_4 = \mathbf{X}_0 + \alpha \mathbf{E}_1 + \beta \mathbf{E}_2 + \gamma \mathbf{E}_3 \quad (3)$$

The  $(\alpha, \beta, \gamma)$  are the coordinates of the point  $\mathbf{X}_4$  in this basis, which are related to the 3-dimensional shape of the object under 3D affine transformations, and are invariant in 3D space. Due to the linear nature of weak perspective  $(\alpha, \beta, \gamma)$  will remain viewpoint invariant under this transformation. We will refer them as 3D

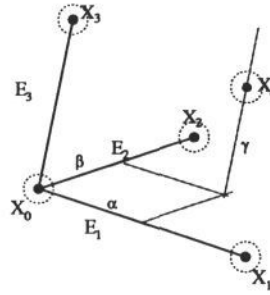


Figure 1: Affine Basis. Each point is surrounded by an error ellipse to show uncertainty in location due to noise and poor localisation

*affine* invariants of our shape. We rewrite (3) in matrix form, and combine it with (2).

$$\begin{bmatrix} su_4 \\ sv_4 \\ s \end{bmatrix} = \begin{bmatrix} t_{11} & t_{12} & t_{13} & t_{14} \\ t_{21} & t_{22} & t_{23} & t_{24} \\ 0 & 0 & 0 & t_{34} \end{bmatrix} \begin{bmatrix} E_{1x} & E_{2x} & E_{3x} & X_{0w} \\ E_{1y} & E_{2y} & E_{3y} & Y_{0w} \\ E_{1z} & E_{2z} & E_{3z} & Z_{0w} \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \gamma \\ 1 \end{bmatrix} \quad (4)$$

Which we can rewrite more simply as:

$$\begin{bmatrix} u_4 - u_0 \\ v_4 - v_0 \end{bmatrix} = \begin{bmatrix} e_{1u} & e_{2u} & e_{3u} \\ e_{1v} & e_{2v} & e_{3v} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix} \quad (5)$$

where  $\mathbf{u}_0, \mathbf{u}_4$ , are the projected basis origin and point respectively,  $\mathbf{e}_1 \dots \mathbf{e}_3$  are the projections of our basis vectors, all of which are obtainable from image measurements. This provides us with two equations and three unknowns, the problem has become underdetermined, so a single 2D image does not allow the recovery of the invariants. A second view with known point correspondences to the first view will however give an overdetermined set of equations solvable by standard least square minimisation methods.

### Solving Overdetermined Equations

Our previous equation (5) can now be written in slightly different matrix form to include the effects of noisy measurements:

$$\begin{bmatrix} u_4 - u_0 \\ v_4 - v_0 \\ u'_4 - u'_0 \\ v'_4 - v'_0 \end{bmatrix} = \begin{bmatrix} e_{1u} & e_{2u} & e_{3u} \\ e_{1v} & e_{2v} & e_{3v} \\ e'_{1u} & e'_{2u} & e'_{3u} \\ e'_{1v} & e'_{2v} & e'_{3v} \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix} + \begin{bmatrix} \vdots \\ \epsilon_i \\ \vdots \end{bmatrix} \quad (6)$$

$e'$  is the projection of the basis in our second view and  $\epsilon$  is the noise in the measured points. Integrating more views into this equation is done trivially by each time adding two more rows to the matrix and can be used to give better conditioned equations. The least squares solution  $\hat{\alpha}$  can then be found using standard methods.

### 3 Object Model Acquisition and Recognition

The object recognition scheme is based on representing an object using these 3D affine invariants. This imposes a slight limitation on the method in that we only consider image data viewed under weak perspective.

Feature extraction as with most recognition schemes forms the first step, in this case we wish to reliably detect corners in the image which relate to the projections of vertices on an object. Multiple views of the object with known point correspondences are then used to calculate the object invariants. In order to avoid the point correspondence and calibration problems of multiple views the vertices are tracked through a sequence of images. A Kalman Filter is then used to optimally estimate the invariants from the stream of data.

The next task is to find invariants from the image data which match those from a model, by having identical values and being associated with the same basis. Instead of calculating the invariants associated with a model at recognition time, they are precomputed and data linked to them is stored off-line, speeding up the final recognition process at the cost of extra storage space.

The geometric hashing method proposed by Wolfson et al [11, 7], with an expanded index to reflect the 3D (as opposed to 2D) object database, is used to implement this data storage. The invariants form a fixed hash table index which points at a *bin* in the table, where we store information about the features and object used to compute those invariants. When attempting to match our image data with a model we extract the contents of the indexed bin, thereby discovering all possible object and pose correspondences. By using several invariants, and hence the contents of several bins, we are able to prune the possible object matches down by finding multiple evidence for particular bases. These can then be examined more closely by trying to match projected model edges with those present in the image.

#### 3.1 Feature detection

In order to simplify object recognition it is often useful to preprocess images, extracting features within the object, which help to give a simplified but distinctive description of it. Corners are useful for describing artificial objects and are of special interest due to their accurate localisation. Sensitivity to noise, speed and selectivity are all issues which need to be considered when attempting to implement a corner detector. A number of methods for detecting corners within shapes include template matching, second-order derivative schemes [12], autocorrelation [13], and median based methods [14].

The corner detector implemented in the following experiments is particularly suited to detecting corners which are projections of vertices. It locates corners by first finding regions where several lines are present and then pinpointing where those lines intersect. Since vertices are the features used in our object description it is helpful to the overall recognition speed, if the corner detector doesn't pinpoint too many other types of corner. The procedure is tolerant to noise and reasonably fast (Figure 2). (For more details see [15].)

#### 3.2 Feature Tracking

A way of avoiding the correspondence problem in multiple images is to exploit temporal continuity in a monocular sequence of frames. We track all the corner points initially detected in the image as the object moves in the scene. Since the corners will only shift slightly between consecutive frames, the approach taken is to perform local corner detection centered on the old corner point locations,

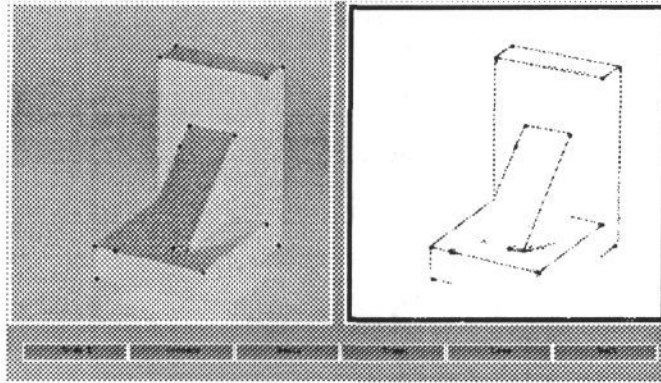


Figure 2: An image of a block with detected corner points superimposed. The image is in reverse video to highlight the corners. There are a couple of noisy points due to dents on the block and other effects such as shadows. Also other points are not detected because our edge threshold was set too high. There is a compromise here in trying to keep background noise to a minimum and detecting all points.

finding new corner locations in the neighbourhood. A further improvement is to perform cross-correlation between a template of the old corner and new corner in cases where we need to disambiguate between several possible maxima in a local region. However it was found that the extra robustness gained by doing this was outweighed by the extra computational cost. Figure 3 shows an example of a test object where corners have been detected and then tracked.

### 3.3 Optimal Estimation of $\alpha, \beta, \gamma$ and their Uncertainty

The most important stage in the recognition scheme is the accurate estimation of the 3D affine invariants. However there is some uncertainty in any estimate of  $\alpha, \beta$  and  $\gamma$  due to errors introduced by the weak perspective approximation and errors in corner localisation. By using the redundancy in a stream of data from multiple views, and ensuring adequate object rotation between the first and final images it is possible to improve the robustness of the resulting estimate over one calculated from a single stereo image pair. A Kalman filter is used to integrate data from multiple views, and optimally estimate  $\alpha, \beta, \gamma$  [15].

The Kalman filter is a recursive linear estimator which successively calculates a minimum variance estimate for a variable, on the basis of observations that are linearly related to this variable. The basic algorithm consists of a continuous "prediction - observation - update" cycle, where the error between prediction and observation is used to revise the old estimate. The initialisation of the filter can either be done either by setting states to zero or by using a batch calculation to work out an initial estimate of the invariants. A covariance matrix representing the initial confidence in the variable estimate is also set.

Our Kalman filter is in fact very simple since the output  $\alpha, \beta, \gamma$  is *invariant*. The Kalman filter will also return the uncertainty of the  $\alpha, \beta, \gamma$  estimates in terms of estimate variances. When this uncertainty falls to within acceptable limits we can stop tracking the corners and proceed with the actual object recognition.

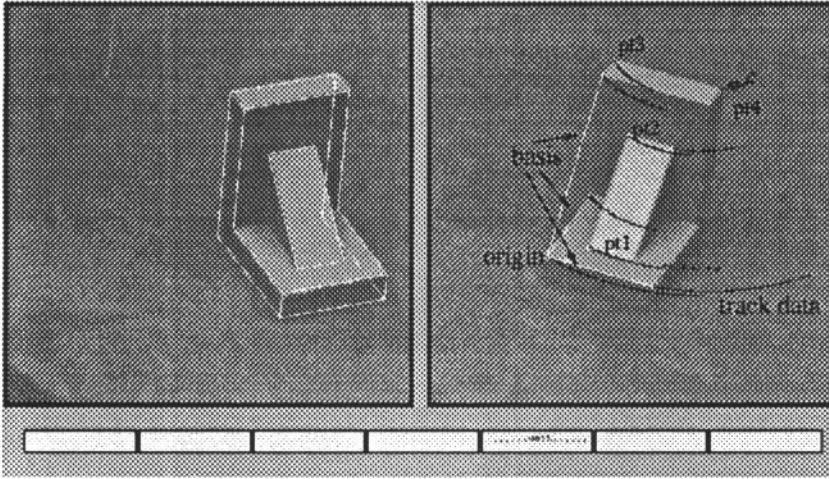


Figure 3: Several points were tracked on our object as it rotated. This is shown as black arcs on our image.

Invariant	Real data	Perfect Result	Invariant	Real data	Perf. Res.
pt1 $\alpha$	0.28	0.28	pt2 $\alpha$	0.29	0.28
$\beta$	0.29	0.29	$\beta$	0.41	0.43
$\gamma$	0.01	0.00	$\gamma$	0.57	0.57
pt3 $\alpha$	0.01	0.00	pt4 $\alpha$	0.98	1.00
$\beta$	0.23	0.27	$\beta$	0.30	0.27
$\gamma$	1.01	1.00	$\gamma$	0.98	1.00

Table 1: Invariants output by the Kalman filter

## Results

Figure 4 shows the results of calculating four sets of invariants using a Kalman Filter (using the track data seen in figure 3). All the graphs have an initial transient stage due to the initialisation of the filter with zeros, however this soon settles down as more track data is incorporated, improving the estimate of the invariants. Points 1 and 4 have the noisier tracks, and this is reflected in the graphs, where the invariant estimate has a greater uncertainty taking longer to settle. The final output is compared with what we expected from perfect data in table 1.

The 3D basis geometry will also affect the sensitivity of  $\alpha, \beta, \gamma$  to noise. If the basis is formed from points which are nearly coplanar the resulting equations (6) which are to be solved will become ill-conditioned. It may be necessary to try several different bases in order to find one which gives well-conditioned equations. This could also be implemented by running several Kalman filters on the track data in parallel, each using a different set of points as the basis and selecting the invariants from the one with least uncertainty.

### 3.4 Geometric Hashing for data storage

The geometric hashing approach to model-based object recognition was proposed by Lamdan, Schwartz and Wolfson [11, 7]. It relies on computing object related

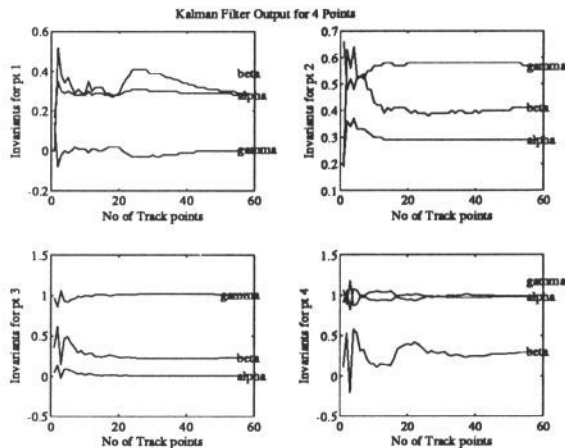


Figure 4: Invariants output by the Kalman filter

invariants from our features, which are then used as a simple and stable index  $H(\alpha, \beta, \gamma)$  into a hash table. Each indexed location in the hash table will store information about the object used to create the invariant, the points used to form the basis and the point referenced in that basis (7).

$$H(\alpha, \beta, \gamma) = (\text{Object}, \text{Basis}, \text{Point}) \quad (7)$$

At the recognition phase we then index into the hash table, and retrieve all possible objects and their bases, using a voting scheme to decide on their likelihood and a final matching process to select the best fitting object.

Initially a large amount of pre-computation is required on the model library to create the hash table, but since this can be done off-line it will not affect recognition speed. Invariants are calculated for all permutations of basis and points and data stored in the hash table. Since 3D data is available in our model we can calculate our invariants directly using (8) [cf (3)].

$$\begin{pmatrix} E_{1x} & E_{2x} & E_{3x} \\ E_{1y} & E_{2y} & E_{3y} \\ E_{1z} & E_{2z} & E_{3z} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} = \begin{pmatrix} X_4 - X_0 \\ Y_4 - Y_0 \\ Z_4 - Z_0 \end{pmatrix} \quad (8)$$

An obvious problem when creating the hash table is the large quantity of data present. If we have an object with  $N_v$  vertices, there are  $N_b = N_v!/(N_v - 4)!$  possible bases and  $N_{invs} = N_b(N_v - 4)$  associated invariant sets. This exponentially growing problem needs to be addressed, by pruning away invariant sets. A number of choices are available:

- Four coplanar points cannot be used to form an adequate basis.
- Ill-conditioned bases will be susceptible to noise and should not be included.
- Use characteristic views of an object [16] to decide which points will be available from a common view, hence removing points never present simultaneously.
- Eliminate any points the corner detector is poor at localising in the real images if this is known.
- An important decrease in data can also be achieved by noting that we do not need to store object information for all invariants calculated from permutations

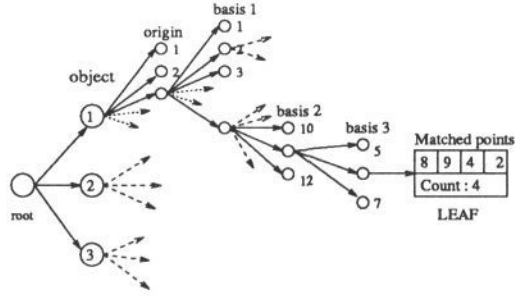


Figure 5: Recognition Tree: For each entry in a bin we register a vote in our tree.

of the same basis. The effect is simply to permute our invariants, so by arranging  $(\alpha, \beta, \gamma)$  in order of size we just need to store information at one location in the hash table. This gives a six-fold reduction in data.

Taking this a step further we note that for a set of four points there are 4! ways of forming a basis. However any invariants calculated in one of these bases is related to all the others [15], so it is possible to find a global invariant for all of these. The drawback with this is we increase the amount of computation that needs to be done at recognition time.

One issue not addressed so far is the effect of poor corner localisation on our invariants  $(\alpha, \beta, \gamma)$ . A study performed by Grimson et al. [17] on 2D affine data shows that if we assume our corner points lie within some error disc, radius  $\epsilon$  (see figure1), the region of uncertainty associated with our affine coordinates is an ellipse whose area, centre and principle axes depend on  $\epsilon$  and the invariants. This would become an ellipsoid in our 3D case. The area of this ellipse increases  $\propto (1 + |\alpha| + |\beta|)^2$ , so larger affine coordinates will have larger errors. This matches well with the worst case study done in [7] which state results from numerical analysis.

$$(A + \delta A) \cdot (\underline{\alpha} + \delta \underline{\alpha}) = (\underline{x} + \delta \underline{x}) \quad (9)$$

Given we can estimate  $\delta A$  and  $\delta \underline{x}$  then  $\delta \underline{\alpha} \leq (1 + x)c_k \epsilon$  (where  $x = |\alpha| + |\beta| + |\gamma|$  is the sum of the absolute values of our invariants, and  $c_k = |c_{k1}| + |c_{k2}| + |c_{k3}|$ , is the sum of the absolute values of the k'th row in the matrix  $C = A^{-1}$ ).

As indicated in the error analysis the error increases proportionally with  $\alpha, \beta, \gamma$ , so we therefore linearly increase the size of the *bins* in our hash table for larger invariants. Bases in our images which are reasonably short are also more affected by poor corner localisation.

### 3.5 The Recognition Process

The actual recognition process is now outlined, and some results are given.

1. Extract corner points from the initial image, these points are then tracked as the object moves.
2. Select four of these points to form the basis, and calculate invariants for all the other points. A good selection of basis will provide a better estimate of the invariants. If the recognition process is a failure a different basis can be tried. Kalman filtering ensures optimal estimation of the 3D invariants.
3. For each invariant set index into the hash table. For every (Object , Basis) grouping which appears in the associated bin register a vote in a recognition tree (figure 5). The recognition tree is shown as a quick way of registering votes for a particular (object, basis) match, which is faster than keeping and



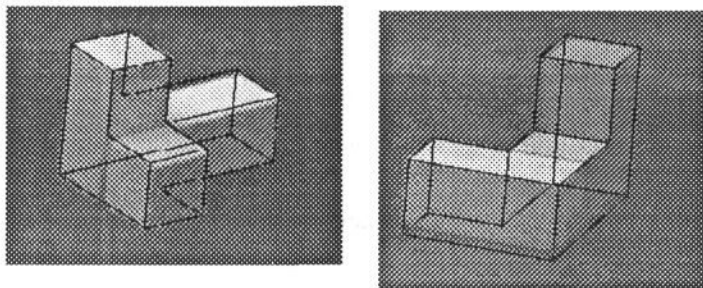


Figure 6: Recognised model is projected onto various shapes.

updating a list of possible matches. Use the object number, origin point and then consecutive basis points to select which branch to follow. At the final leaf register a vote for this grouping by incrementing a counter. The point match for this invariant set is also stored here. In order to account for uncertainty in the 3D invariant estimate, examine not just the indexed bin but also local surrounding ones.

4. If a grouping (Object , Basis) scores a large number of votes, then this object is possibly present in the scene.
5. Compute the transformation (2) between matched model and real data points. A minimum of four points are required to calculate this transformation. At least five point matches (usually more) are available, so least squares minimisation is used to find a solution. Verify that the least squares solution to the transformation projects all model points satisfactorily to their matched image points.
6. Transform the edges of the model according to this transformation. Verify how close the projected edges lie to scene edges, if there is good correlation then the object has been recognised. Models that have been successfully projected onto the image are shown in figures 3 and 6. Failure to match the edges requires the process to be restarted with a different basis.

## 4 Summary and Conclusions

An object recognition scheme has been proposed and implemented which is successful at identifying simple 3D objects. Central to the system is an object representation that is viewpoint invariant under the weak perspective assumption, allowing fast matching with a model database. The representation is inherently robust to partial occlusion, since the invariants are obtained from groups of just five corner points. By tracking more features we provide enough information to determine a small number of possible object matches, which can then be evaluated by matching image edges with projected model edges. Certain issues remain which include increasing the speed and reliability of the tracking system to ensure our 3D invariants are as accurate as possible. We also wish to further examine the problem of exponential growth of data in the hash table, which occurs with increasing numbers of features on an object.

The aim of our future work is test the performance of this system with more complex objects and in cluttered scenes. We will be evaluating the discriminating power and computational efficiency of the method under such conditions.

## References

- [1] R.T.Chin and C.R.Dyer. Model-based recognition in robot vision. *ACM Computing Surveys*, 18(1):67-108, 1986.
- [2] P.J.Besl and R.C.Jain. Three-dimensional object recognition. *ACM Computing Surveys*, 17(1):75-145, 1985.
- [3] D.G.Lowe. The viewpoint consistency constraint. *Int. Journal of Computer Vision*, 1:57-72, 1987.
- [4] D.W. Thompson and J.L. Mundy. Three-dimensional model matching from an unconstrained viewpoint. In *Proceedings of IEEE Conference on Robotics and Automation*, pages 208-220, 1987.
- [5] W.E.L.Grimson and T.Lozano-Perez. Localising Overlapping Parts by Searching the Interpretation Tree. *IEEE Trans. Pattern Analysis and Machine Intell.*, 9(4):469-482, 1987.
- [6] J.L. Mundy and A.Zissermann editors. *Geometric Invariance in Computer Vision*. MIT Press, 1992.
- [7] Y.Lamdan, J.T.Schwartz, and H.J.Wolfson. Affine invariant model-based object recognition. *IEEE Trans. on Robotics and Automation*, 6(5):578-589, 1990.
- [8] L.G.Roberts. Machine perception of three - dimensional solids. In J.T. Tippet, editor, *Optical and Electro-optical Information Processing*. MIT Press, 1965.
- [9] J.J.Koenderink and A.J.van Doorn. Affine structure from motion. *J. Opt. Soc. America*, 8(2):377-385, 1991.
- [10] E.B.Barrett et al. In J.L. Mundy and A.Zissermann, editors, *Geometric Invariance in Computer Vision*, chapter 14, pages 277-292. MIT Press, 1992.
- [11] H.J.Wolfson and Y.Lamdan. In J.L. Mundy and A.Zissermann, editors, *Geometric Invariance in Computer Vision*, chapter 17, pages 335-353. MIT Press, 1992.
- [12] L.Kitchen and A.Rosenfeld. Gray-level corner detection. *Pattern Recog. Letters Vol.1*, pages 95-102, 1982.
- [13] C.G. Harris. In A. Blake and A. Yuille, editors, *Active Vision*, chapter 16. MIT Press, 1992.
- [14] K.Paler et al. Local ordered grey levels as an aid to corner detection. *Pattern Recognition Vol.17 No.5*, pages 535-543, 1984.
- [15] S. Vinther and R. Cipolla. Towards 3D object model acquisition and recognition using 3D affine invariants. Technical Report CUED / F - INFENG / TR136, Dept. of Engineering, University of Cambridge, 1993.
- [16] H.Plantinga and C.R.Dyer. Visibility, occlusion and the aspect graph. *Int. Journal of Computer Vision*, 5(2):355-39, 1990.
- [17] D.P.Huttenlocher W.E.Grimson and D.Jacobs. A study of affine matching with bounded sensor error. In *2nd European Conf. on Computer Vision*, 1992.