# Model-Based 3-D Object Recognition Using Scalar Transform Descriptors.

Kenneth M. Dawson and David Vernon.

Trinity College, University of Dublin, Department of Computer Science,
Dublin 2, Ireland.

## ABSTRACT

Three dimensional object recognition is an essential capability for any advanced machine vision system. We present a new technique for the recognition of 3-D objects on the basis of comparisons between 3-D models. Secondary representations of the models, which may be considered as complex scalar transform descriptors, are employed. The use of these representations overcomes the common dependency of matching individual model primitives (such as edges or surfaces). The secondary representations used are one-dimensional histograms of components of the visible orientations, depth maps and needle diagrams. Matching is achieved using template matching and normalised correlation techniques between the secondary representations. We demonstrate the power of this new technique with several examples of object recognition of models derived from actively sensed range data.

## 2. BACKGROUND

Object recognition is the process of identifying material things by associating what is sensed with what is known. This association is typically performed by a hypothesis generation and verification procedure, in which objects which are known *a priori* are compared to sensed, or *viewed*, objects.

For any general purpose object recognition strategy, with the exception of a few approaches which consider every possible view of every *known* object as a separate model[1,2,3], it is typically accepted that *known* object models must be three-dimensional. However there are two distinct approaches to the consideration of *viewed* objects; one in which a 2-D (image based) model is employed and the other in which a 3-D model is extracted, either using the various depth-from-... routines or using actively sensed range data. While it is possible to recognise objects using 2-D models[4], recognition is (or at least, intuitively, should be) facilitated through the use of 3-D models. It is on the basis of this premise that we addressed the second of these approaches and present, in this paper, a technique for the recognition of rigid 3-D structures in 3-D *viewed* data.

The derivation of 3-D models from 3-D *viewed* data is typically regarded as a problem of segmentation - into physically relevant surfaces. Object recognition, then, is usually addressed as a graph-matching task in which connected graphs of surfaces are compared. However most researchers who have addressed the problems of surface fitting and segmentation have found them to be far from trivial[5,6,7,8,9]. For example, Hoffman *et al.*[10], developed a system which classifies surfaces in range data as planar, convex or concave, and justified this simple classification on the basis that more complex classifications are extremely sensitive to noise.

This problem of the stability of segmentation and classification in three-dimensional data does not occur in the case of the Extended Gaussian Image (or EGI)[11,12,13], as the EGI represents all of data from a three-dimensional model. However, the EGI does not uniquely identify objects and matching of EGIs has proved difficult[14,15]. Drawing though, on the concept behind the EGI and on other work of Horn[16] and Ikeuchi[17], we employ several different representations of an entire model to roughly approximate, tune and verify hypotheses about *viewed* data. This new approach overcomes the problems of the EGI and provides a reliable measure for hypothesis verification.

Each of the representations used is comprised of an array of scalars and so, in some ways, may be regarded as scalar transform descriptors. In addition, as entire models are compared in single operations rather than explicitly matching model primitives the method, described herein, has been dubbed *implicit model matching*.

# 3. IMPLICIT MODEL MATCHING

## 3.1 Introduction

The basic problem of object recognition is the correspondence of an object model which is *viewed* with a particular *known* object model. This operation typically also involves computing the associated pose of the corresponding *known* object model, a problem which has six degrees of freedom; three wrt. orientation and three wrt. position. The technique of implicit model matching employs various secondary representations in order to consider this problem in terms of sub-problems.

Initially, potential *known* object orientations are approximated by comparing each possible view of a given *known* object, as defined by a tessellated sphere, in terms of the visible orientations, with those visible from the *viewed* object. Potential orientations are then fine tuned by correlating 1-D histograms of orientations (known as directional histograms). Object position is approximated using simple geometry, and fine tuned using a template matching technique, and finally each hypothesis is evaluated/verified by comparing needle diagrams generated from the relevant views of both the *viewed* and the *known* models.

At each stage in the generation, tuning and verification of hypotheses only comparisons of the various secondary representations are employed. The central concept behind this technique, of implicit model matching, is then that 3-D object models may be reliably compared through the use of secondary representations, rather than through the use/comparison of model component primitives (e.g. surfaces).

## 3.2 Object Domain and Scene Complexity

The object domain being addressed is that of rigid objects, and (initially at least) only objects which have been segmented from all other objects in the scene are considered. The *viewed* scene data is only constrained to provide a means of generating the various secondary representations (and, of course, of allowing objects to be segmented from all others in the scene). The secondary representations, which are directional histograms, needle diagrams and depth maps, may be generated from any view of any 3-D model. Hence it is possible to test the approach using a view of a theoretical object model, or equally using a 3-D model derived from active (or passive) range data.

## 3.3 Approximating Object Orientation

The first stage in this technique is the approximate determination of orientations of a *known* object model which may, potentially, correspond to the *viewed* object model. This is achieved by considering the *known* object from every possible viewpoint, as defined by a tessellated sphere, and comparing directional histograms of tilt (which will be explained presently) for every viewpoint with a directional histogram of tilt derived from the *viewed* model. The orientations which generate locally maximum correlations between these histograms (i.e. as compared to the correlations associated with neighbouring tessellations) may be regarded as the potentially matching orientations.

### 3.3.1 Directional Histograms

The concept of the Directional Histogram was developed as part of the technique of implicit model matching and embodies the notion of mapping a single component of the 3-D orientations of a model visible from a given viewpoint to a 1-D histogram, where the component of orientation is defined about the axes of the viewing device. Four different components may be employed: roll, pitch, yaw and tilt; where roll, pitch and yaw are defined as rotations about the Z, X and Y axes of the viewing device respectively, and tilt is defined as $\pi$ less the angle between the orientation vector and the focal axis (i.e. the Z axis of the viewing device); See Figure 1. Example directional histograms are shown in Figures 2 and 3 in order to demonstrate the usefulness and the practicalities of this representation.
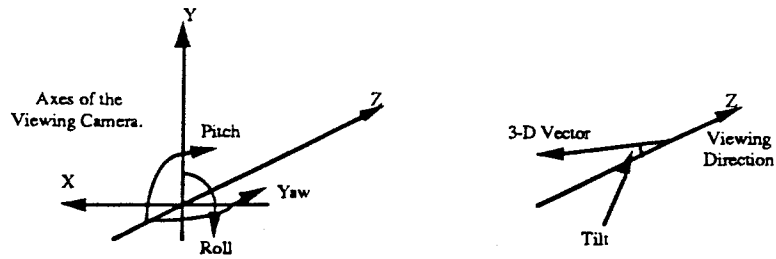
**Fig. 1.** Definitions of roll, pitch, yaw and tilt about the coordinate axes of the viewing device.
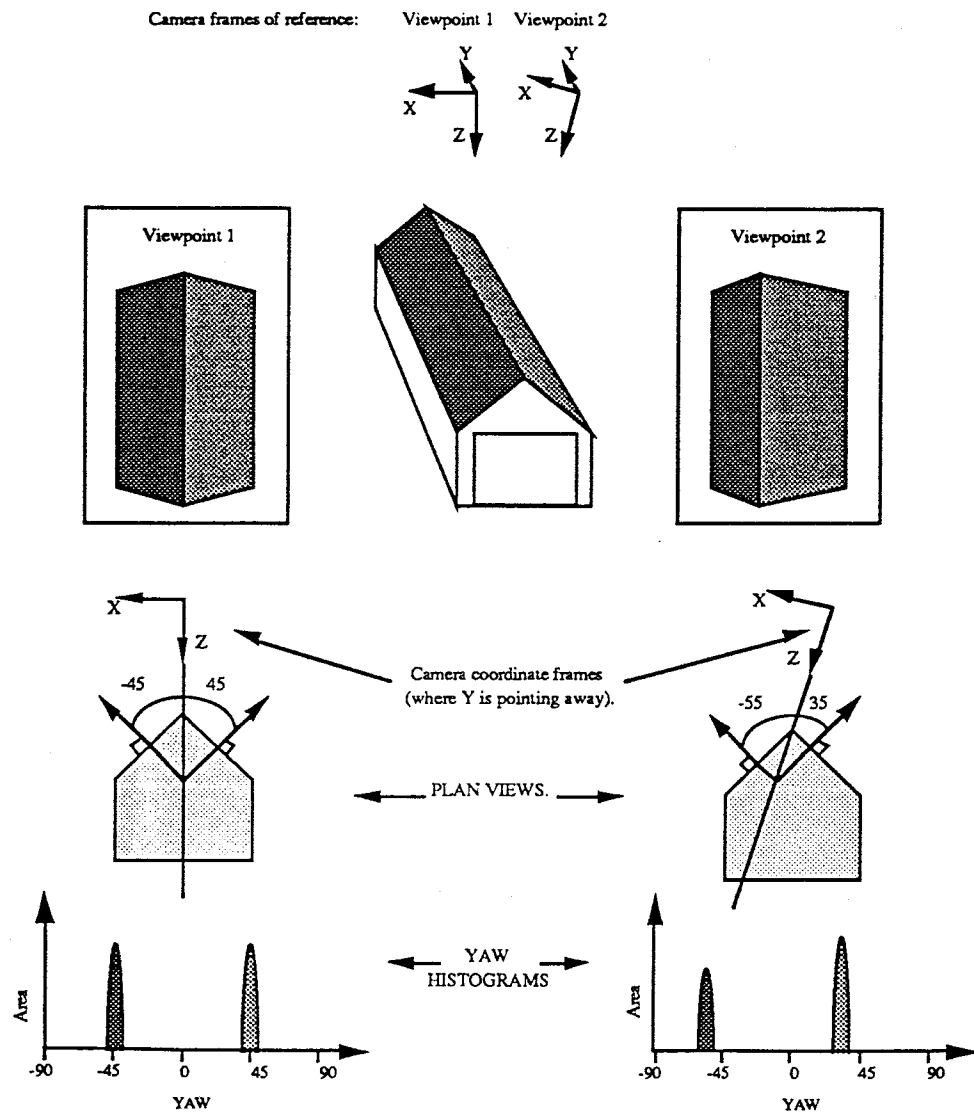


**Fig. 2.** Example yaw directional histograms. These two yaw histograms of a garage-like object are a simple example of how directional histograms work. The visible surface areas of the views of the object are mapped to the histograms at their respective yaw angles (which are defined with respect to the axes of the viewing device/camera). In order to aid understanding the sections of the histograms which are mapped from given surfaces are shaded in the grey-level of the surfaces. Notice the slight shift in the histograms for the two different viewpoints.
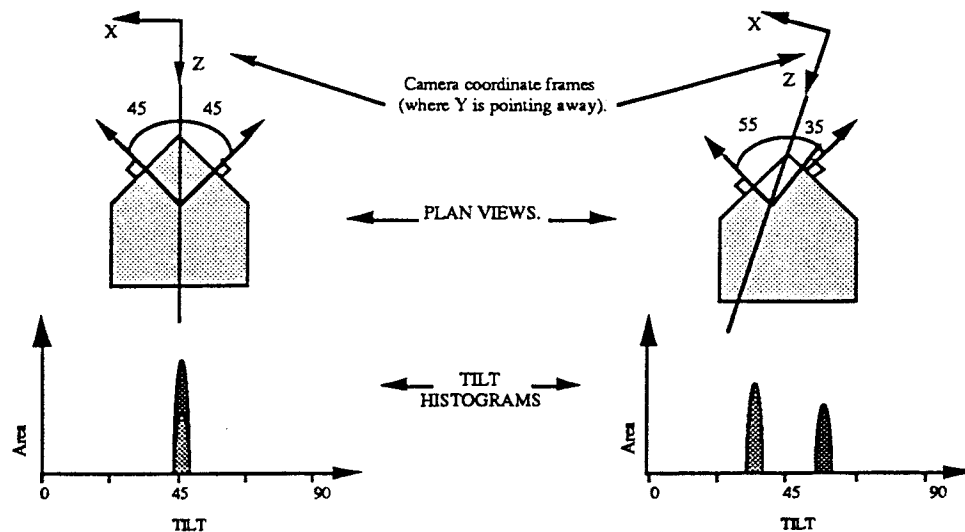
Fig. 3. Example tilt directional histograms. These two tilt histograms are derived from the two views of the garage-like object shown in Fig. 2. Notice how, for the first view the two orientations result in the same value of tilt, and in the second view how the values change.

### 3.3.2 Resulting Orientations

The result of these comparisons of tilt directional histograms is the identification potentially matching orientations. However, only two degrees of freedom have been constrained, as only the tilt of the object is approximated. In order to complete the approximation of orientation, we must also compute potentially matching values of roll around the focal axis of the viewing device; See Figure 4.
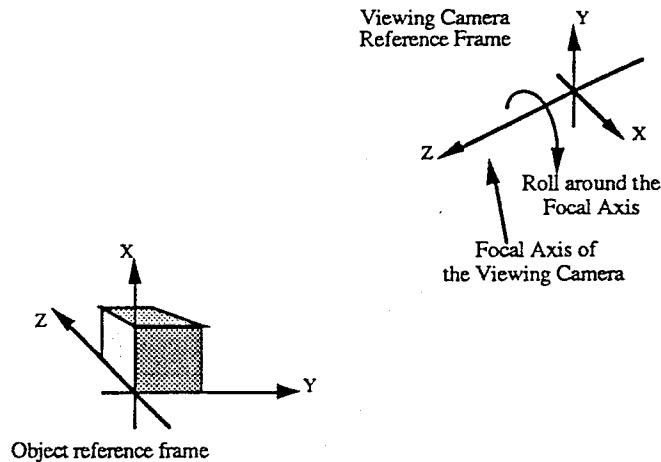


Fig. 4. Viewing angle geometry.

This identification may again be performed using directional histograms, but using the roll component of orientation, rather than that of tilt. For every determined value of tilt for a *known* model, a directional histogram of roll is derived and compared (using normalised cross correlation) with that from the *viewed* model, in every possible value of roll (i.e. by simply shifting the *known* models' roll histogram in a circular fashion, and comparing with the *viewed* models' roll histogram after each shift). The local maxima correlations once again identify the potentially matching rolls/orientations.

## 3.4 Fine-tuning Object Orientation

The potentially matching orientations computed can only be guarantied to be as accurate as the quantisation of the sampled sphere. Increasing the resolution of the sphere to an arbitrarily high level, however, would obviously cause a significant increase in the computational overhead required in determining potential orientations. Alternatively, though, it is possible to fine-tune the orientations using directional histograms (of roll, pitch and yaw) in a similar fashion to the method used for the approximate determination of object roll. Pitch, yaw and roll directional histograms are derived from the view of the *known* object and compared with pitch, yaw and roll directional histograms derived from the *viewed* model. The differences between the directional histograms indicate the amount by which the orientation may best be tuned (e.g. see Figure 2). The various directional histograms are derived and compared both sequentially and iteratively until the tuning required falls below the required accuracy of orientation or until the total tuning on any component of orientation exceeds the range allowed (which is defined by the quantisation of the tessellated sphere).

This stage allows the accuracy of potentially matching orientations to be determined to an arbitrarily high level (limited only by the resolution of the directional histograms - which can be made arbitrarily high). Hence, although only a limited number of possible viewpoints of any *known* object are considered, the orientation of the object may be determined to a high level of accuracy.

## 3.5 Approximating Object Position

Turning now to the approximate determination of object position (relative to the viewing device/camera), it is relatively straightforward to employ the position of the *viewed* model with respect to its viewing camera. The imaged centroid of the *viewed* model, and an approximate measure of the distance of the *viewed* model from the viewing camera are both easily computed. The position of the camera which views the *known* model may then be approximated by placing the camera in a position relative to the known model's 3-D centroid, such that the centroid is at the correct approximate distance from the camera and is viewed by the camera in the same position as the *viewed* model's imaged centroid. See Figure 5.
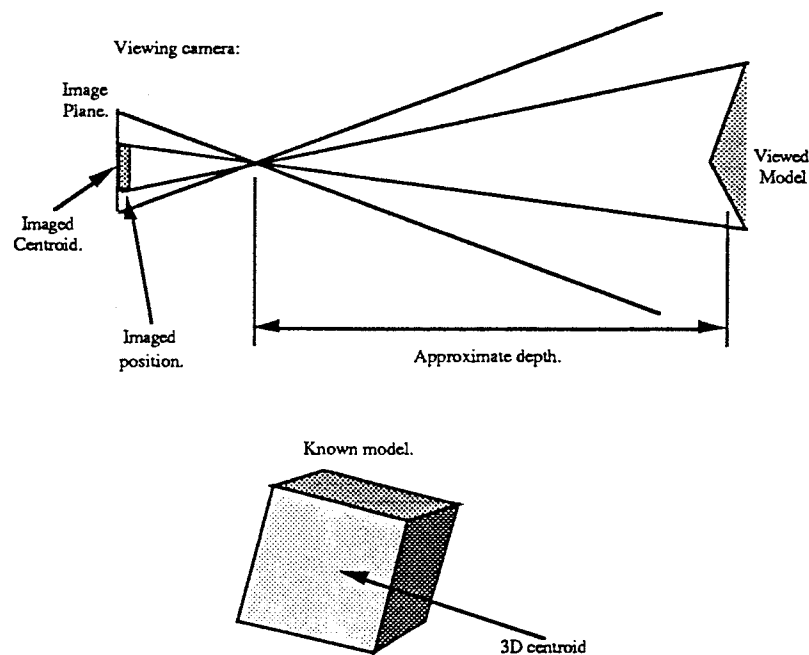


Fig. 5. The geometrical features used when determining approximate object position.

## 3.6 Fine-tuning Object Position

Fine tuning object position may be considered in terms of two operations; tuning position in a directional orthogonal to the focal axis of the viewing device, and tuning of the distance of the object from the same viewing device (i.e. the depth). This separates the 3 degrees of freedom inherent in the determination of object position.

### 3.6.1 Tuning Object viewed position (i.e. orthogonal to viewing device)

This operation is performed using a template matching technique in which a needle diagram (i.e. an iconic representation of the visible local surface orientations) of the *known* model is compared using a normalised correlation mechanism, with a needle diagram of the *viewed* model. The position of the template which returns the highest correlation is taken to be the optimal position for the *known* model.

The standard method of comparing iconic representations is normalised cross correlation, but this form of correlation is defined only for scalars. For the comparison of needle diagrams 3-D vectors must be compared and that correlation (*NV*) for each possible position of the template *(m,n)* is defined as follows:

$$NV(m,n) = \frac{\sum_i \sum_j f(viewed(i,j)) * (\pi - angle(viewed(i,j), known(i-m,j-n)))}{\sum_i \sum_j f(viewed(i,j)) * \pi}$$

where *viewed(i,j)* and *known(i,j)* are the 3-D orientation vectors from the *viewed* and *known* needle diagram respectively, *f(vector) = 1* if the vector is not NULL (and 0 otherwise), and *angle(vector1,vector2)* is the angle between the two vectors. This definition results in a maximum of one if the two needle diagrams are identical, and a minimum of zero if they have no needles in common positions.

In order to make this template matching operation more efficient, the needle diagrams are first compared at lower resolutions. These comparisons use only a simple measure *S(m,n)* of the number of needles in the *viewed* models' needle diagram which have a corresponding needle (or any orientation) in the *known* models' needle diagram:

$$S(m,n) = \frac{\sum_i \sum_j f(viewed(i,j)) * f(known(i-m,j-n))}{\sum_i \sum_j f(viewed(i,j)) * \pi}$$

where the same definitions apply. Only positions of the template for which *S(m,n)* is greater that 0.75 are considered at higher resolutions using *NV*.

### 3.6.2 Tuning Object Depth

Fine tuning the distance between the *known* model and its viewing camera is done through direct comparison of the depth maps generated from both the *viewed* model and the *known* model (in its determined pose). The Depth Change Required (or *DCR*) is defined as follows:

$$DCR = \frac{\sum_i \sum_j f(viewed(i,j)) * f(known(i,j)) * (viewed(i,j) - known(i-m,j-n))}{\sum_i \sum_j f(viewed(i,j)) * f(known(i,j))}$$

where *f(depth) = 1* if the depth is defined and 0 otherwise. This depth change is directly applied as a translation to the pose of the camera which views the *known* model, in a direction defined by the focal axis of the camera. Due to perspective effects this operation will have effects on the depth map rendered, and so is applied iteratively until the *DCR* falls below an acceptable level.

## 3.6.3 Final tuning of Object Position

Tuning of object position is limited in accuracy primarily by the resolution used in the needle diagrams (for tuning position orthogonal to the focal axis of the viewing camera). As a final stage in tuning we attempt to overcome this by determining the position to sub-pixel accuracy. This is accomplished using a combination of template matching, normalised correlation and quadratic modelling techniques. The best position may be determined to pixel accuracy using the technique described in section 3.6.1. In order to determine the position to sub-pixel accuracy the normalised correlations around the best position (as determined to pixel accuracy) are used and are modelled as quadratics in two orthogonal directions (i.e. parallel to the two image axes). See Figure 6.
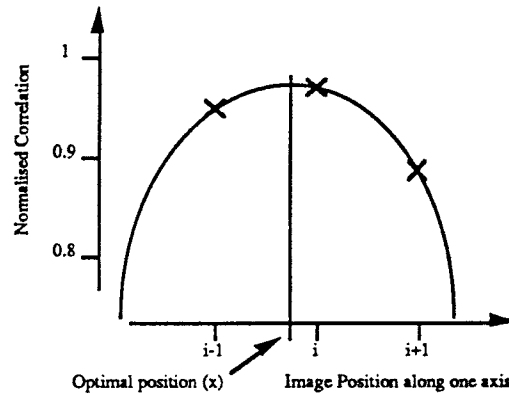


Fig. 6. Modelling correlations around a central maximum using quadratic modelling.

The quadratic equation used is $y = a.x^2 + b.x + c$ where $y$ is the normalised correlation determined with the template in position $x$ along one image axis, relative to the best position $i$ determined to pixel accuracy along that axis. Considering the best position and the positions on either side:

$$
\begin{array}{ccccccc}
x_{i-1}^2 & x_{i-1} & 1 & & a & & y_{i-1} \\
x_i^2 & x_i & 1 & * & b & = & y_i \\
x_{i+1}^2 & x_{i+1} & 1 & & c & & y_{i+1} \\
X & & & * & C & = & Y
\end{array}
$$

Applying the pseudo inverse, C may be calculated:
$$ C = (X^T.X)^{-1}.X^T.Y $$
Now, taking the derivative of $y$ with respect to $x$, we get:
$$ dy/dx = 2.a.x + b $$
and $y$ will be maximised when $dy/dx$ is 0, so the value of $x$ which maximises the normalised correlation is:
$$ x = -b / (2.a) $$

## 3.7 Verifying Hypotheses

Having hypothesised and fine tuned poses of *known* objects it is necessary to determine some measure of fit for each hypothesis so that it may be accepted (subject to no better hypothesis being determined), or rejected. The normalised correlation of local surface orientations (i.e. needle diagrams) between the *viewed* model and the *known* model in a determined pose, as used when fine tuning object position (See section 3.6.1) gives a degree-of-fit which represents all aspects of object position and orientation. This degree-of-fit, then, provides a powerful hypothesis verification measure.

The pose of the *known* model associated with the best degree-of-fit, given that it the degree-of-fit is over a high threshold is taken to represent the *viewed* model. Search for this best degree-of-fit is continued until all invoked possibilities are considered.

## 3.8 An exception

There is, however, one situation in which this technique, of implicit model matching, will fail and that is when only a single surface/orientation is visible. Determination of object roll in this instance is impossible using directional histograms (as there is an inherent ambiguity with respect to roll around the orientation vector of the surface).

Detecting this situation may be done by considering the standard deviation of the visible orientations as mapped to an Extended Gaussian Image (as if the s.d. is less than a small angle then it may be taken that only one surface is visible).

Coping with this situation is quite straightforward, as when only one orientation is visible the problem may be regarded as one of shape recognition. The only problems caused by this situation are the approximation and tuning of object roll, and it is possible to employ signatures of boundary curvature in place of the directional histograms of roll. These may be used in much the same fashion as the directional histograms (See Figure 7), and the rest of the algorithm functions as before.
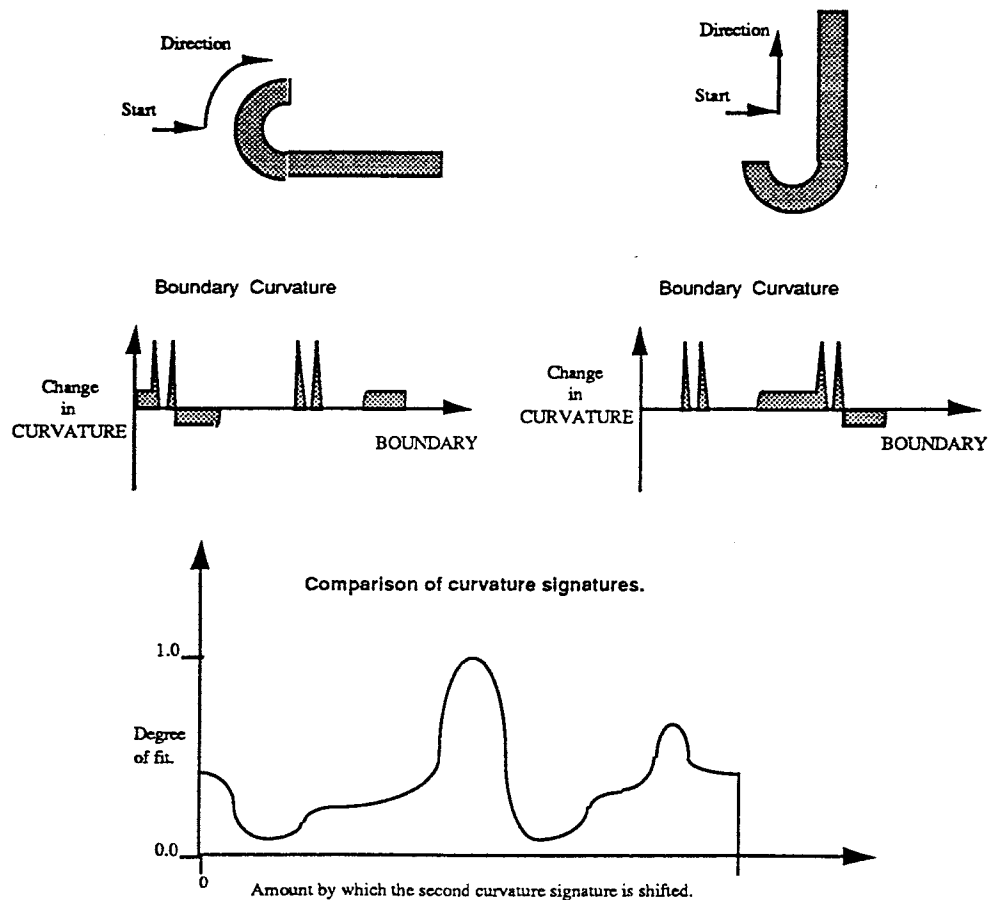
Fig. 7. Example comparison of boundary curvature.

# 4. EXPERIMENTAL RESULTS

The basic technique of implicit model matching was presented in the previous section. Reiterating, it is a technique in which 3-D object models are compared using various different secondary representations. It was also stated that the technique can be tested with any form of three dimensional model from which the secondary representations can be generated. Initial testing was performed using views of theoretical objects, and subsequently was extended to testing with models derived from actively sensed range data. The use of actively sensed range data was justified on the basis that it provides reasonably accurate 3-D models and results in "an excellent domain to study issues without trivializing the problem and without unnecessary complexity" as Ramesh Jain points out in the preface of Ting-Jung Fan's book[14]. Some of the results of that (second set of) testing are presented here. It should be noted that all of the actively sensed range data was supplied by the Pattern Recognition and Image Processing Lab. of Michigan State University.

Details of how three-dimensional models are extracted from the range data (e.g. segmentation techniques, etc.) are not addressed in this paper, and the interested reader is directed towards [18] or [19]. The known models are specified using a simple CAD specification technique (again see [18] or [19]) and are shown in Figure 8.
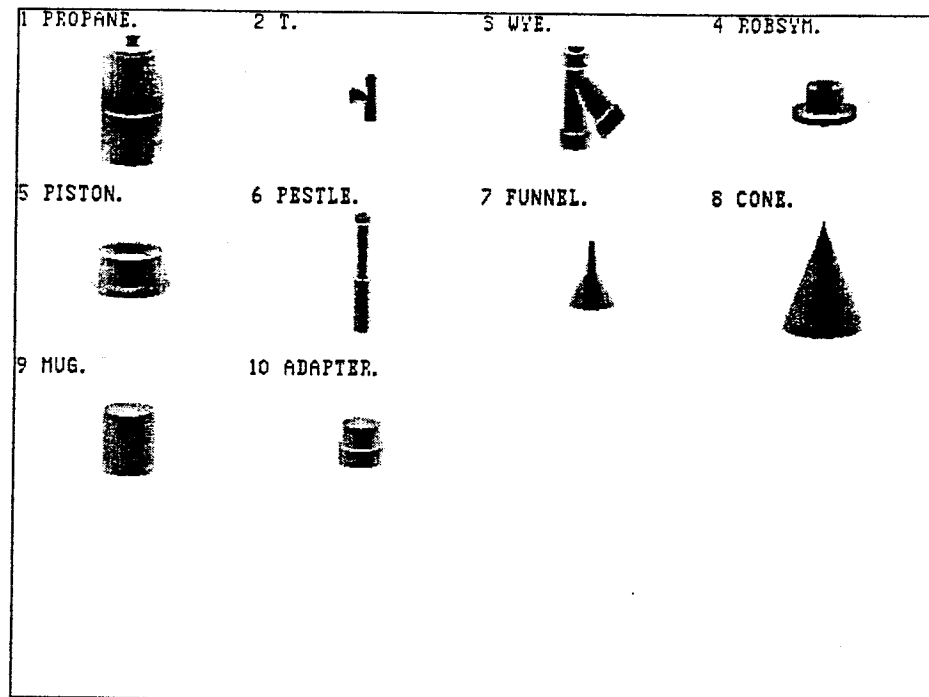


Fig. 8. CAD Models of the known objects.

Finally in Figures 9, 10 and 11 examples of recognised objects are shown. In each figure a rendering of the needle diagram and of the various directional histograms of the *viewed* model (i.e. derived from the actively sensed range data) are shown, in the top left and right quadrants respectively. In the bottom left quadrant, a rendering of the needle diagram of the best matching *known* model in its determined pose is shown, and in the bottom right quadrant a subtraction of the rendered needle diagrams is given.

THE VIEW TO BE RECOGNISED:       TILT.           ROLL.

PITCH.        YAW.

POTENTIAL MATCHES:
BEST MATCH FOUND:   FUNNEL.
(DEG. = 0.9637).

BEST MATCH FOUND:   FUNNEL.
(DEG. = 0.9637).
NEEDLE SUBTRACTION:

Fig. 8. Recognition of a funnel.

THE VIEW TO BE RECOGNISED:       TILT.           ROLL.

PITCH.        YAW.

POTENTIAL MATCHES:
BEST MATCH FOUND:   ADAPTER.
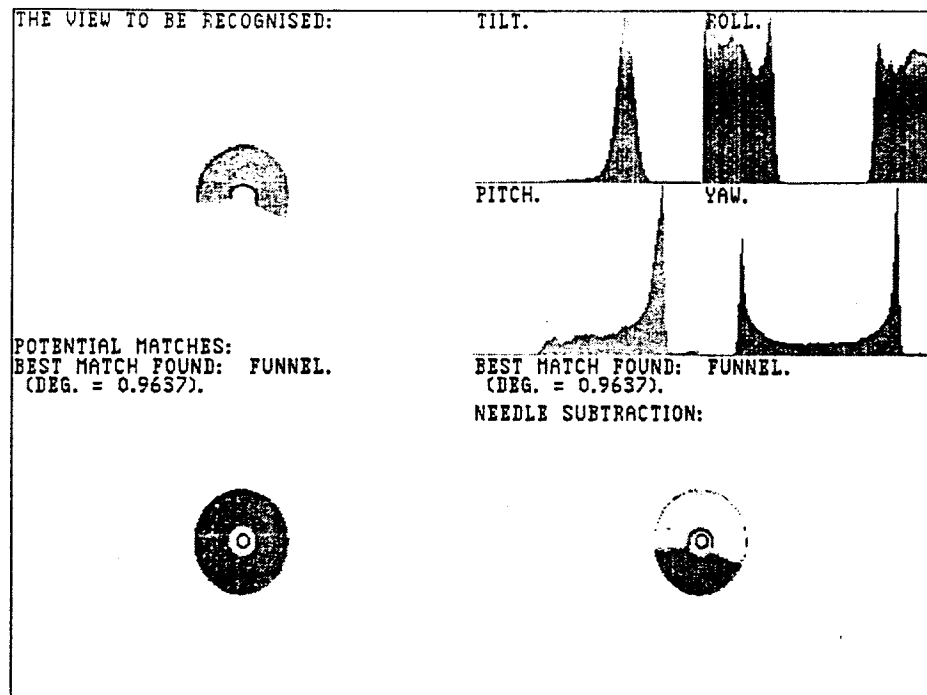(DEG. = 0.9503).

BEST MATCH FOUND:   ADAPTER.
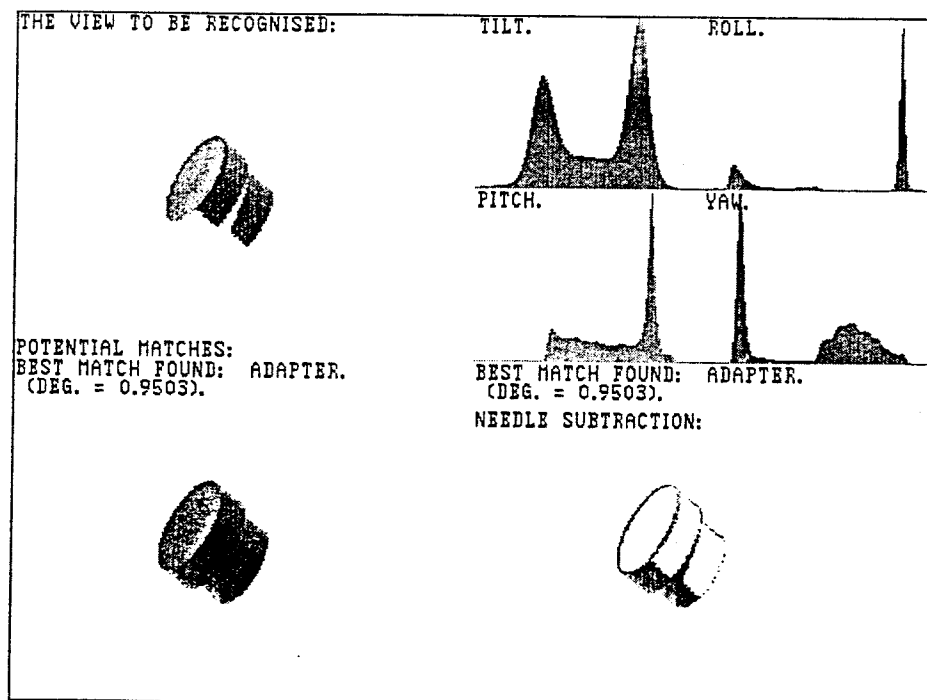(DEG. = 0.9503).
NEEDLE SUBTRACTION:
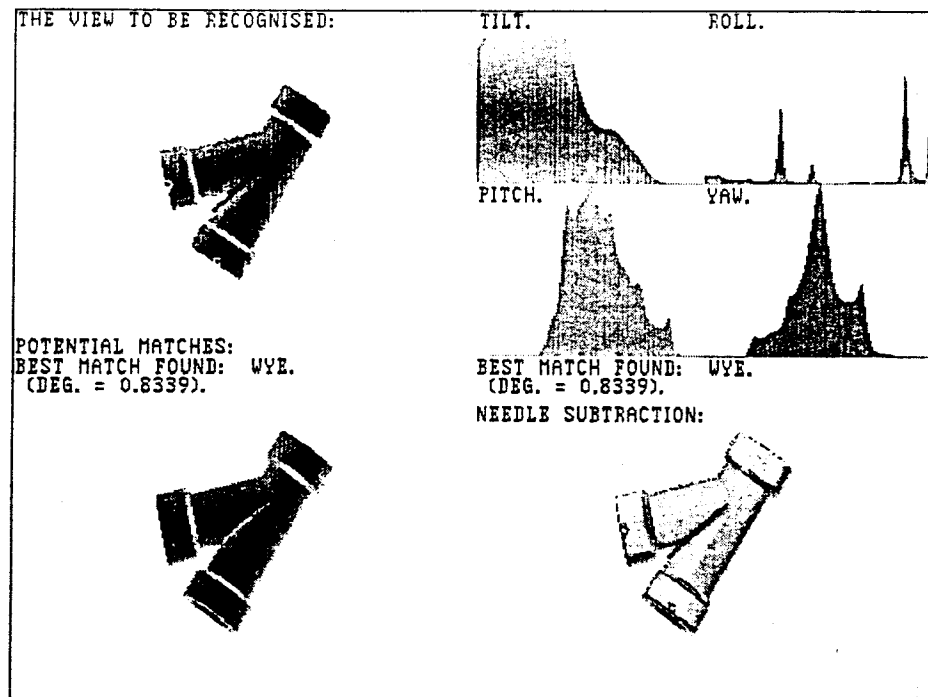
Fig. 9. Recognition of an adapter.

Fig. 10. Recognition of a wye piping joint.

## 5. CONCLUSIONS

The examples shown demonstrate the ability of the technique, of implicit model matching, to correctly identify an object and its pose for single objects. However it is also important to address the discriminability of the approach (i.e. how well each object is discriminated from the other objects, and for this reason a table of the best degrees-of-fit between each *viewed* instance and each *known* model was drawn up and is shown in Figure 11.

| Scene | Figure | Known models and best associated degrees-of-fit when compared with the scene model. | | | | | | | | | | Recognised model |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Piston | Mug | Cone | Funnel | Robsym | Adapter | Propane | Wye | Pestle | T | |
| Piston | | 0.8238 | 0.5836 | 0.6246 | 0.2692 | 0.4934 | 0.4565 | 0.7230 | 0.4357 | 0.1752 | 0.1630 | Piston ✓ |
| Coffee mug | | 0.5719 | 0.9340 | 0.8087 | 0.3009 | 0.5514 | 0.5046 | 0.8805 | 0.5470 | 0.1754 | 0.1975 | Mug ✓ |
| Cone | | 0.6040 | 0.6466 | 0.9492 | 0.4749 | 0.3870 | 0.2827 | 0.8220 | 0.5953 | 0.2840 | 0.1503 | Cone ✓ |
| Funnel | | 0.6772 | 0.5785 | 0.7149 | 0.8864 | 0.6190 | 0.5551 | 0.6275 | 0.7640 | 0.6980 | 0.3809 | Funnel ✓ |
| Funnel | 8 | 0.7962 | 0.9340 | 0.8264 | 0.9637 | 0.7599 | 0.7898 | 0.8260 | 0.7772 | 0.3030 | 0.3746 | Funnel ✓ |
| Robsym | | 0.5825 | 0.5877 | 0.5413 | 0.3351 | 0.8759 | 0.5464 | 0.6077 | 0.5071 | 0.2157 | 0.2077 | Robsym ✓ |
| Adapter | 9 | 0.8185 | 0.8592 | 0.7917 | 0.5981 | 0.7798 | 0.9503 | 0.8973 | 0.5961 | 0.2596 | 0.2589 | Adapter ✓ |
| Adapter | | 0.6855 | 0.8813 | 0.8170 | 0.6940 | 0.6905 | 0.9069 | 0.8073 | 0.6831 | 0.2859 | 0.3561 | Adapter ✓ |
| Propane cylinder | | 0.4509 | 0.5019 | 0.5932 | 0.1912 | 0.3549 | 0.2389 | 0.9378 | 0.4251 | 0.1788 | 0.0979 | Propane ✓ |
| Wye piping joint | 10 | 0.5102 | 0.3992 | 0.5639 | 0.1894 | 0.4040 | 0.2620 | 0.5085 | 0.8339 | 0.1581 | 0.1506 | Wye ✓ |
| Wye piping joint | | 0.5527 | 0.3196 | 0.5579 | 0.1867 | 0.3387 | 0.1830 | 0.5811 | 0.7634 | 0.0792 | 0.0291 | Wye ✓ |
| Pestle | | 0.4710 | 0.3892 | 0.6092 | 0.3088 | 0.3510 | 0.2392 | 0.6831 | 0.6458 | 0.8422 | 0.2614 | Pestle ✓ |

Fig. 11. Table of results of the best degrees of fit.

Each viewed instance was correctly identified, although as can be seen from figure 11, in a number of instances high measures were associated with several *known* models. For example, consider the 0.8805 degree-of-fit between the propane *known* and the view of the coffee mug. This is caused by the fact that part of the propane model (i.e. part of the cylinder) is similar to the view of the coffee mug which was considered, and also by the definition of *NV* (the correlation of needle diagrams) which takes no account of the regions of the *known* model which do not correspond to the *viewed* model. This second cause could be resolved by changing the formula for *NV*, but it should be realised that the reasoning behind the current definition was to allow for possible occlusions.

In conclusion, the technique of implicit model matching is demonstrated to be quite robust in the domain of models derived from actively sensed range data, confirming the hypothesis that 3-D objects may be reliably recognised through the comparison of secondary representations.

## 6. REFERENCES

1. T-J. Fan, <u>Describing and Recognising 3-D Objects using Surface descriptors</u>, Springer-Verlag, New York, 1990.

2. B. Bhanu, ``Representation and Shape Matching of 3-D Objects," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 6, No. 3, pp. 340-351, May 1984.

3. W.E.L Grimson, ``On the Recognition of Curved Objects," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 11, No. 6 pp. 632-643, June 1989.

4. D.G. Lowe, ``Two-Dimensional Object Recognition from Single Two-Dimensional Images," *Artificial Intelligence*, Vol. 31, No. 3, pp. 355-395, March 1987.

5. R.B. Fisher, ``Using Surfaces to recognise partially obscured objects," *Proceedings of the 8th International Joint Conference on Pattern Recognition, Karlsruhe, West Germany*, pp. 989-995, 8-12 August 1983.

6. R.B.Fisher, <u>From Surfaces to Objects</u>, Wiley, Chicester, 1989.

7. P.J. Flynn and A.K. Jain, ``Surface Classification: Hypothesis Testing and Parameter Estimation," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Ann Arbor*, pp. 261-267, 1988.

8. O. Faugeras, M. Herbert and M. Pauchon, ``Segmentation of range data into planar and quadratic patches," *Proceedings of the Conference on Computer Vision and Pattern Recognition, Washington*, pp. 8-13, 1983.

9. P. Besl, <u>Surfaces in Range Image Understanding</u>, Springer-Verlag, New York, 1988.

10. R. Hoffman and A.K. Jain, ``Segmentation and classification of range images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 9, No. 5, pp. 608-620, September 1987.

11. B.K.P. Horn and K. Ikeuchi, ``Picking Parts out of a Bin," *MIT Artificial Intelligence Laboratory*, Memo no. 746, October 1983.

12. B.K.P. Horn, ``Extended Gaussian Images," *Proceeding of the IEEE*, Vol. 72, No. 12, pp. 1671-1686, December 1984.

13. K. Ikeuchi, ``Determining Attitude of Object from Needle Map using Extended Gaussian Image," *MIT Artificial Intelligence Laboratory*, Memo no. 714, April 1983.

14. K. Ikeuchi, N.K. Nishihara, B.K.P Horn, P. Sobalvarro, and S. Nagata, `` Determining Grasp configurations using photometric stereo and the PRISM binocular stereo system," *The international Journal of Robotics Research*, Vol. 5, No. 1, pp. 46-65, Spring 1986.

15. P. Brou, ``Using the Gaussian Image to Find Orientation of Objects," *The International Journal of Robotics Research*, Vol. 3, No. 4, pp. 849-125, Winter 1984.

16. B.K.P. Horn and B.L. Bachmann, ``Registering Real Images using Synthetic Images," <u>Artificial Intelligence - An MIT Perspective</u>, pp. 129-160, 1979.

17. K. Ikeuchi, ``Generating an interpretation tree from a CAD model for 3D-object recognition in bin-picking tasks," *International Journal of Computer Vision*, Vol. 1, No. 2, pp. 145-165, 1987.

18. K. Dawson, ``Three-Dimensional Object Recognition through Implicit Model Matching," Ph.D. thesis, Dept. of Computer Science, University of Dublin, Trinity College, Dublin, Ireland, 1991.

19. K. Dawson, ``3-D Representations" and ``3-D Object Recognition," <u>Parallel Computer Vision: The *VIS à VIS* System</u>, D. Vernon and G. Sandini (editors), Chapters 8 and 9, Ellis Horwood (to be published).