

HUANG Ying, DING Xiao-qing, WANG Sheng-jin

## Recognition of 3-D objects based on Markov random field models

© Higher Education Press and Springer-Verlag 2006

**Abstract** The recognition of 3-D objects is quite a difficult task for computer vision systems. This paper presents a new object framework, which utilizes densely sampled grids with different resolutions to represent the local information of the input image. A Markov random field model is then created to model the geometric distribution of the object key nodes. Flexible matching, which aims to find the accurate correspondence map between the key points of two images, is performed by combining the local similarities and the geometric relations together using the highest confidence first method. Afterwards, a global similarity is calculated for object recognition. Experimental results on Coil-100 object database, which consists of 7 200 images of 100 objects, are presented. When the numbers of templates vary from 4, 8, 18 to 36 for each object, and the remaining images compose the test sets, the object recognition rates are 95.75 %, 99.30 %, 100.0 % and 100.0 %, respectively. The excellent recognition performance is much better than those of the other cited references, which indicates that our approach is well-suited for appearance-based object recognition.

**Keywords** Pattern recognition, 3-D object recognition, Markov random field, Highest confidence first

### 1 Introduction

In real-world scenes, the representation of a 3-D object may be modified due to multiple factors such as: 1) object scale, viewpoint and illumination variations; 2) partial occlusion; 3) noise data; 4) object deformation. Humans can distinguish different objects easily without considering

these variations. However, this is quite a difficult task for computer vision systems. Most object recognition approaches aim to find object features and then match these features between the observed data and the object databases. Generally, this problem can be traced back to establish the relation between two images. In the following we will present a novel method to solve this visual correspondence problem.

There are numerous research efforts dealing with the object recognition problem and the existing approaches can mainly be classified into two distinct categories: approaches based on global feature and approaches based on local feature. The approaches based on global feature extract global image features such as color histograms [1], and receptive field histograms [2]. Global features are robust to scale and viewpoint changes, but it has been difficult to extend them to partially occluded and noise images. Methods based on support vector machines [3–4], Sparse Network of Winnows (SnoW)[9], and eigenspace matching [5] can handle images corrupted by noise and partial occlusions successfully, but fail to recognize viewpoint variant or heavily occluded objects. The approaches based on local feature identify local feature points, and then create a local image descriptor at each interest point [6–7]. Object matching is performed by finding similar point pairs between test images and training models. Advantages of the local image features are that they are only partially affected by object occlusion, viewpoint modification, and deformation. However, their recognition results depend on the accuracy of the point detection. Points detected in one object may be missed in another image of the same object. Furthermore, no geometric constraint between interest points is utilized to rectify the matching results.

Our method uses dense local features that sampled at a large number of repeatable locations to represent objects. Markov random field models are established to model the geometric constraints between object key points. The matching program is composed of two procedures: local matching and global matching. Local matching procedure calculates the similarities between the key points of two images. In the global matching procedure, the highest confidence first method is introduced to reach a local

---

Translated from *Journal of Tsinghua University (Science and Technology)*, 2005, 45(1):28-32 (in Chinese)

---

HUANG Ying, DING Xiao-qing(✉), WANG Sheng-jin  
Department of Electronic Engineering, Tsinghua University, Beijing  
100084, China  
E-mail: dxq@ocrserv.ee.tsinghua.edu.cn

minimum of the MRF model. This final result decides which pairs of key points correspond to the same point and which points have no corresponding partner in the other image. Our object recognition method is evaluated using the Coil-100 object database containing 7 200 images of 100 objects. Different numbers of images are selected as template examples, and the remaining as test images. We also test our method on the images corrupted by synthetically generated rotation, scaling and occlusions. The remarkable recognition rates show the potential of our approach in the problem of recognizing 3-D objects.

The remainder of the paper is organized as follows: the following section presents our approach in detail. Sect.3 describes the local and global matching procedures. Some experiments and comparisons are given in Sect.4. We conclude this paper in Sect. 5.

## 2 General framework

### 2.1 Motivation

As mentioned above, objects can be represented by local features. These local features are computed at some interest points, which can be extracted using corner detectors. Object matching is performed by comparing the features of the key points between two images. However, the detected positions of these key points will be influenced by image transformations. Thus for two images of the same objects, some points in the first image cannot find their accurate partners in the second one. In addition, the geometric positions of the key points have some relations even after complex transformations. For example, considering a key point in the first image, the accurate partners of its neighbor key points are also located in the neighborhood of its partner in the second image. This information should be considered.

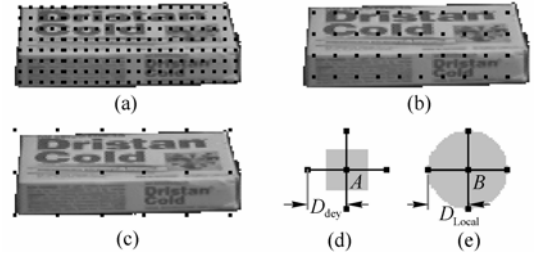
To solve these problems, this paper introduces two kinds of topological grids. The first kind of grid is utilized to extract local image features, while the other grid combines these local features with the geometric constraints between grid nodes. Unlike previous approaches, our method does not identify key locations. Thus we do not have to consider the accuracy of point detection, and the algorithm complexity is decreased. In our framework, the interest points are equidistantly distributed on different levels, and these points are sampled densely enough to cover most of the object details.

### 2.2 MRF-based framework

We call the first kind of grid “local grid”, and the second kind of grid “key grid”. The nodes of these grids are equidistantly sampled. The local grid is composed of several levels of local nodes (Fig. 1(a), (b)). The distance between two adjacent local nodes of a sampling level  $k$  is

defined as  $D_{\text{local}}(k)$  (Fig. 1(e)). The key grid is introduced to perform the match between two images (Fig. 1(c)).  $D_{\text{key}}$  is the distance between two adjacent key nodes (Fig. 1(d)):

$$D_{\text{key}} = \frac{F(\text{image})}{V_{\text{stad}}} \quad (1)$$



**Fig. 1** (a), (b): Two level local grid; (c) Key grid; (d) Neighbourhood of the key node  $A$  (gray regions); (e) Neighborhood of the local node  $B$  (gray regions)

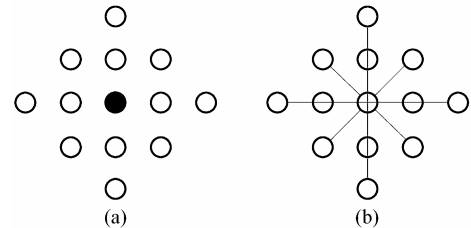
$V_{\text{stad}}$  is a standard value that decides the resolution of the key grid, and  $F(\text{image})$  is calculated from some physical parameters of the input image, such as the image pixel number, or the image bounding box.  $D_{\text{local}}(k)$  is formulated as:

$$D_{\text{local}}(k) = \frac{0.5D_{\text{key}}k}{K}, k = 1, 2, \dots, K-1, K \quad (2)$$

where  $K$  is the maximum level number. For example,  $K$  is equal to 2 in Fig. 1.

Since  $D_{\text{key}}$  is much larger than  $D_{\text{local}}(k)$ , in the neighborhood of each key node (Fig. 1(d)), there are several local nodes. These local nodes are used to calculate the similarity between two key nodes in the matching procedure.

Given two images of the same objects, assume that the neighbourhood of two key nodes  $A$  and  $B$  correspond to the same region on the physical object. The corresponding partners of the neighbour key nodes of  $A$  should also be adjacent to the node  $B$ . On the other hand, if we have already determined the partners of the neighbour nodes of  $A$ , then the coordinate of  $B$  depends mainly on two factors: one is the similarity between  $A$  and  $B$ , and the other is the positions of these partners. Therefore, this issue can be modeled by a Markov random field model. As shown in Fig. 2, a 12-node neighbourhood and its associated clique are defined in our MRF model. More details of the potential calculation and the relaxation algorithm will be described in Sect. 3.2.



**Fig. 2** (a) 12-Node neighbourhood; (b) Its associated clique

### 2.3 Color histograms

Researchers have developed many local image descriptions, such as color histograms, Gabor features, differential invariants, and scale invariant features. To accelerate the test experiments, we take use of the color histogram description.

Color histograms are popularly used in many applications because they are trivial to compute, and robustly tolerate image transformations and changes in camera viewpoint. We quantize colors into a set of  $L$  representative colors  $\mathcal{C} = \{c_1, c_2, \dots, c_L\}$ . For each local node, we compute the color histogram feature using its neighbor pixels.

Two color histograms  $H_1$  and  $H_2$  are compared by computing their intersection, an idea introduced by Swain and Ballard [1]. The intersection is

$$I(H_1, H_2) = \sum_{i=1}^L \min[H_1(i), H_2(i)] \quad (3)$$

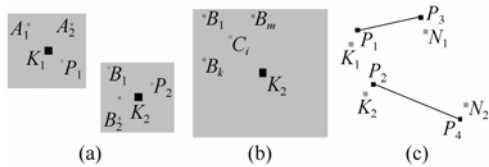
This intersection is equal to a value between 0 and 1. We define this data as the similarity  $S_{\text{local}}$  between two local nodes.

## 3 Template matching

### 3.1 Local matching

The objective of the local matching procedure is to find the similarity  $S_{\text{key}}$  between two key nodes. This data is calculated from the similarities between the local nodes that is located in the neighborhood of the two key nodes.

In Fig. 3(a), two key nodes  $K_1$  and  $K_2$  are displayed. Local nodes located in their neighborhood are signed as  $(A_1, A_2, \dots, A_N)$  and  $(B_1, B_2, \dots, B_M)$ , respectively.  $S_{\text{local}}(i, j)$  is defined as the similarity between the local node  $A_i$  and  $B_j$ . To obtain  $S_{\text{key}}(K_1, K_2)$ , we first extract data  $S_{A_i}$ ,  $i = 1, 2, \dots, N$  for  $A_1, A_2, \dots, A_N$ , which describes the similarities between these nodes and their optimum matching points in  $K_2$ 's neighborhood. Since local nodes are uniformly and densely distributed, the largest value of  $S_{\text{local}}(i, j)$ ,  $j = 1, 2, \dots, M$  can be selected as  $S_{A_i}$ . The corresponding local node, denoted as  $B_i$ , is set as the local matching point of  $A_i$ .



**Fig. 3** Local and global matching (see text for details). (a), (b): Matching between two key nodes  $K_1$  and  $K_2$ ; (c) Geometric potential calculation

Then we use these calculated data to obtain the similarity  $S_{\text{key}}(K_1, K_2)$  between the two key nodes.  $K_1$ 's accurate

matching point may not be  $K_2$  because key nodes of the two images are sparsely sampled. We should find two best matching points in the neighborhood of  $K_1$  and  $K_2$ , which correspond to the same physical point on the object. Therefore, we combined the local matching point  $B_i$  of  $A_i$  ( $i = 1, 2, \dots, N$ ) with the local similarities  $S_{A_i}$  to compute the two best matching points  $P_1$  and  $P_2$  (Fig. 3(a)):

$$X_{P_1} = \frac{\sum_{i=1}^N S_{A_i} X_{A_i}}{\sum_{i=1}^N S_{A_i}}, \quad X_{P_2} = \frac{\sum_{i=1}^N S_{A_i} X_{B_i}}{\sum_{i=1}^N S_{A_i}} \quad (4)$$

To compute the local similarity  $S_{\text{key}}(K_1, K_2)$ , all the neighbor pixels of  $K_1$  are considered. Each point may be located in the neighbor regions of multiple local nodes. Among the similarity of these nodes, we can obtain a maximum data. Then  $S_{\text{key}}$  is formulated as:

$$S_{\text{key}}(K_1, K_2) = \frac{\sum_{D \in N(K_1)} \max_{N(A_i) \ni D, 1 \leq i \leq N} (S_{A_i})}{N_{\text{point}}} \quad (5)$$

Here  $D$  is a neighbor pixel,  $N(K_1)$  denotes the neighbor points of  $K_1$ ,  $N(A_i)$  denotes the neighbor points of  $A_i$ .  $N_{\text{point}}$  is the number of pixels in  $N(K_1)$ .

### 3.2 Global matching

Given two images  $I_1$  and  $I_2$ , the local matching procedure has calculated the similarity between every two key nodes of the two images. The global matching procedure combines these results with  $I_1$ 's MRF model to decide the corresponding map between the key nodes of  $I_1$  and  $I_2$ .

Our MRF model is composed of:

- A set of sites  $S = \{s_1, s_2, \dots, s_n\}$ . Each site corresponds to a key node of  $I_1$ .
- A set of possible labels  $\Omega = \{l_1, l_2, \dots, l_m\}$ . These labels  $l_1, l_2, \dots, l_m$  represent the  $m$  key nodes of the second image  $I_2$ .
- A neighborhood relation over the sites, which defines a graph where the vertices represent sites.

Let  $\omega_i$ ,  $\omega_i \in \Omega$ , denotes the corresponding label of the site  $s_i$ . Then we have the conditional posterior potential for each site [8]:

$$E_i(\omega_i | s_i) = V(\omega_i | s_i) + \sum_{c:i \in c} V_c(s_i | s_{N_i}) \quad (6)$$

where  $c:i \in c$  means any clique  $c$  containing the site  $s_i$ . Assume that  $\omega_i = l_j$ , then  $V(\omega_i | s_i)$  is the negative of the similarity between the key nodes  $s_i$  and  $l_j$  of the two images. The other item gives the geometric potential of the current node. The item contains only one component because only one clique is defined (Fig.2(b)).  $s_{N_i}$  defines the corresponding labels of the 12 neighbor nodes of the site  $s_i$ .

Assume that in Fig. 3(c), the key node  $K_1$  corresponds to

the site  $s_i$ ,  $K_2$  corresponds to the label  $l_j$ .  $P_1$  and  $P_2$  are their best matching points.  $N_1$  is a neighbor node of  $K_1$ . Its accurate partner is  $N_2$ .  $P_3$  and  $P_4$  are their matching points. Two parameters can be obtained from the lengths and directions of the two lines  $P_1P_3$  and  $P_2P_4$ :

$$r(P_1P_3, P_2P_4) = \frac{\text{Len}(P_1P_3)D_{\text{key}}(K_2)}{\text{Len}(P_2P_4)D_{\text{key}}(K_1)} \quad (7)$$

$$\text{ddir}(P_1P_3, P_2P_4) = \text{dir}(P_1P_3) - \text{dir}(P_2P_4)$$

where  $D_{\text{key}}$  is defined by Eq. (1). The two data describe the length ratio and direction difference between the two lines. If  $K_2$  is  $K_1$ 's accurate partner, for each one of the 12 neighbor nodes of  $K_1$ , the two parameters should remain approximately unchanged. Furthermore, the length ratio is close to 1.0. Hence the potential  $V_c(s_i | s_{N_i})$  is formulated as:

$$V_c(s_i | s_{N_i}) = -\exp[-\text{var\_dangle} \times w_1 - \text{var\_r} \times w_2 - (\text{avg\_r} - 1)(\text{avg\_r} - 1)w_3] \quad (8)$$

here  $w_1$ ,  $w_2$  and  $w_3$  are three weights,  $\text{avg\_r}$  is the average of the 12 ratios,  $\text{var\_r}$  is their variance, and  $\text{var\_dangle}$  is the variance of the twelve direction differences.

The highest confidence first (HCF) algorithm [8] is a deterministic relaxation technique which can be converged to a local minimum of the MRF, but induce drastically less computational cost than a stochastic relaxation scheme. The HCF algorithm classifies the sites into two classes: "committed" and "uncommitted". Initially all sites are set uncommitted, and a site has no effect on its neighbors unless it has been committed. Thus, in Eqs. (7) and (8), only committed sites are used to compute  $V_c$ .

A stability measure is calculated for each site based on the local conditional posterior potential defined in Eq. (6). This measure determines the order in which the sites are to be visited. Details of this algorithm have been described in Ref. [8]. The procedure terminates when the criterion function can no longer be decreased by reassignment of the labels.

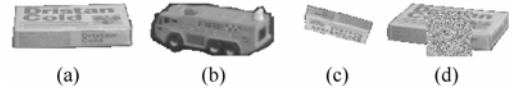
After global matching, we obtain a posterior potential for each key node. The negative of the potential describes the local and geometric similarity between the current node and its partner. Therefore, the average value of these data is used to define a global similarity for the following object recognition experiments.

## 4 Experiments

Our object recognition algorithm is tested on the columbia object image library(COIL) [5] database. The COIL database consists of 7 200 images of 100 objects (72 views for each of the 100 objects). The images are color images of  $128 \times 128$  pixels. The objects are positioned in the center of a turntable and observed from a fixed viewpoint. For each object, the turntable is rotated of  $5^\circ$  per image (Fig. 4(a), (b)).

Parts of the COIL images are selected as the template images. The other images are set as the test set. Tests on corrupted images are also taken. For each test image, we

match it with all the template images. The object ID of the template image with the largest global similarity is returned as the object recognition result.



**Fig. 4** Coil-100 object database. (a), (b): Two images; (c) A synthesized image under scaling and rotation; (d) A synthesized image occluded by a randomly placed  $k \times k$  window of uniformly distributed random noise ( $k = 48$ )

We test the proposed system using all the 100 objects. Different numbers of views are selected as the template set. For each object, the numbers of template images vary from 4 (one every  $90^\circ$ ), 8 (one every  $45^\circ$ ), 18 (one every  $20^\circ$ ) to 36 (one every  $10^\circ$ ). The remaining images compose the test sets. The image pixel number is selected to calculate  $F(\text{image})$  in Eq. (1). In addition, since the resolution of the key grid is decided by the standard value  $V_{\text{stad}}$ , performance of the framework with different  $V_{\text{stad}}$  is tested. Comparisons between our method and other approaches are listed in Tables 1 and 2.

**Table 1** Recognition rates using different template sets/ %

Template images	400	800	1 800	3 600
Test images	6 800	6 400	5 400	3 600
Our framework (color & $V_{\text{stad}} = 25$ )	95.75	99.30	100.0	100.0
Our framework (color & $V_{\text{stad}} = 9$ )	93.66	97.89	99.41	99.98
SNoW (edges)[9]	88.28	89.23	94.13	96.25
SNoW (intensity)[9]	81.46	85.13	92.31	95.81
Nearest neighbor (color)	79.74	93.08	98.91	99.89
Linear SVM (color) [3]	83.99	95.36	99.31	100.0

**Table 2** Comparisons between our framework and other methods[4] (Template images: 400)/ %

Representation	Our framework ( $V_{\text{stad}} = 25$ )	Columbia [5]	Roobaert [4]
Color only	95.75	77.5	82.3
Shape & color	---	87.6	86.9

In Table 1, the image representations of the associated methods are displayed in the parentheses. Results of SNoW are cited from [9]. Results of the methods in Table 2 are cited from [4]. For the nearest neighbor classifier and support vector machines [3], we reduce the image spatial resolution from  $128 \times 128$  to  $32 \times 32$  and then transform each COIL image into an eight-bit vector of  $32 \times 32 \times 3 = 3 072$  components.

Our method performs recognition with excellent percentages of success even in the presence of very similar objects. The

ability of handling viewpoint changing is much higher than the other cited methods. The average match time between two images is no more than 30 ms for the low-resolution framework ( $V_{\text{stad}} = 9$ ), and 200 ms for the high-resolution framework ( $V_{\text{stad}} = 25$ ) on a P4 1.7 GHz computer. However, it still took us quite a long time to finish the test of the high-resolution framework. Therefore, the low-resolution framework is adopted to continue the following test, and eight images for each object are selected as the template set.

Our method is invariant to image shifting, scaling and rotation. This ability is evaluated using synthetically generated images (Fig. 4(c)). The 7 200 synthesized images are tested on the template set, and the recognition rate is 97.94 %.

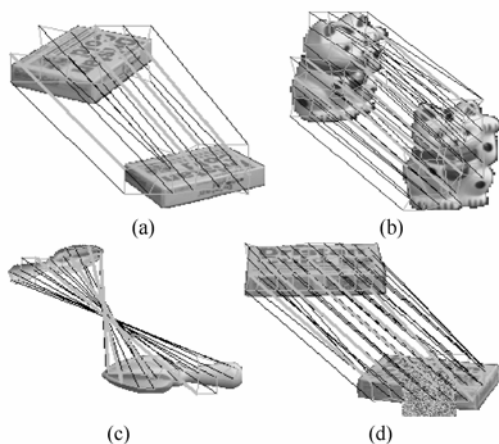
In order to verify the systems against occlusion, test images corrupted by generated occlusions are also synthesized (Fig. 4(d)). Table 3 lists the results of our method and the SVM-based method. From the obtained experimental results, we conclude that our method achieves very good rates even if half pixels of the images are occluded.

**Table 3** Recognition rates for COIL images occluded by a randomly placed  $k \times k$  window of uniformly distributed random noise/ %

$k$	24	48
Our method	97.93	96.72
Linear SVM [3]	95.32	84.99

A shortcoming of our object recognition algorithm is the ability against noise corruption. However, this drawback is mainly caused by the color histograms. Gabor features, which are robust to additive Gaussian noise, will be used to further enhance our system.

Finally we show some correspondence maps between the key nodes of different COIL images in Fig.5.



**Fig. 5** Key node correspondence maps (in dark lines) between different images of the same objects. The matching points of two adjacent key nodes are linked by gray lines. Some mapping lines are emphasized (a), (c)  $V_{\text{stad}} = 9$ ; (b), (d):  $V_{\text{stad}} = 25$

## 5 Conclusions

The main contribution of this paper is to introduce a general framework, which can combine the local image descriptions with the geometric structure of an object together to establish point correspondence maps between different images. Markov random field is introduced to model this framework. Object recognition is then performed from these maps. This method is successfully tested on the Coil-100 image database. The remarkable experimental results indicate that this approach is well-suited for 3-D object recognition robust under viewpoints changing, occlusion, rotation, and scaling. More efforts will be done to further extend this method to detect objects in complex environments.

## References

1. Swain M. J., Ballard D. H, Color indexing. In IJCV, 1991, 7(10):11–32
2. Schiele B., Crowley J. L, Object recognition using multidimensional receptive field histograms, In Proc. ECCV, 1996
3. Pontil M., Verri A, Support vector machines for 3D object recognition, In IEEE Trans. Pattern Analysis and Machine Intelligence, 1998, 20(6):637–646
4. Roobaert D., Van Hulle M, View-based 3D object recognition with support vector machines, In IEEE international workshop on neural networks for signal processing, 1999
5. Murase, Hiroshi, Nayar S. K, Visual learning and recognition of 3-D objects from appearance, In IJCV, 1995, 14 (1): 5–24
6. Lowe D. G. Object recognition from local scale-invariant features, In Proc. ICCV, 1999
7. Jugessur D., Dudek G, Local appearance for robust object recognition, In Proc. CVPR, 2000
8. Chou P. B., Brown C. M, The theory and practice of bayesian image labeling, In IJCV, 1990, 4(3): 185–210
9. Yang M. H., Roth D., Ahuja N, Learning to Recognize 3D Objects with SnoW, In Proc. ECCV, 2000