

# Resolving Occlusion in Augmented Reality : a Contour Based Approach without 3D Reconstruction

M.-O. Berger  
CRIN/CNRS & INRIA Lorraine  
BP 239, 54506 Vandoeuvre les Nancy, France  
e-mail:berger@loria.fr

## Abstract

*We present a new approach for resolving occlusions in augmented reality. The main interest is that it does not require 3D reconstruction of the considered scene. Our idea is to use a contour based approach and to label each contour point as being "behind" or "in front of", depending on whether it is in front of or behind the virtual object. This labeling step only requires that the contours can be tracked from frame to frame. A proximity graph is then built in order to group the contours that belong to the same occluding object. Finally, we use some kind of active contours to accurately recover the mask of the occluding object.*

## 1 Introduction

Augmented reality systems aim at enhancing the user's vision with computer generated images [7, 6, 2]. Most of the time, such systems simply overlay computer generated images and only attempt to minimize object registration errors. However, such methods are effective only when there are no occlusions between real objects and computer generated objects. If a real object occludes a virtual object  $V$ , we must compute the occlusion mask  $m_V$ , that is the region of the image plane to which the visible part of  $V$  corresponds. The virtual object is then only displayed on this mask. Most of the time, occlusions are resolved interactively [3, 2] by allowing the user to delineate the occlusion mask. Hence, solving automatically the occlusion problem for augmented reality is a challenge, especially when little is known about the world to be augmented.

**Context:** We are concerned with video sequences of static scenes taken by a mobile camera (or a mobile observer). Satisfying image composition requires a good temporal registration between the real scene and the virtual objects, as the camera moves. This can basically be achieved using two approaches. The first solution is to use position sensors; but instrumenting the real world is not always pos-

sible. The other solution is to use modeled objects to perform registration.

As we do not want to modify or to instrument the environment with which we interact, we use model based registration. This implies that the 3D model of some features in the scene is known. In fact, this is not too restrictive: for numerous applications, the 3D structure of some objects in the scene are often known (ground, main objects . . . ). The transformation between the camera and the scene frame can then be computed using these features. Let us now get back to the occlusion problem itself.

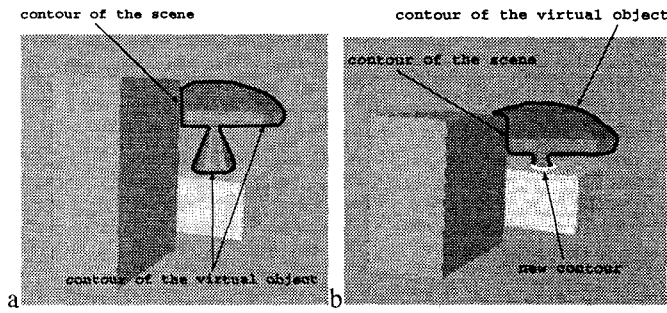
**Related works:** While several augmented reality systems are described in the literature, few of them address the occlusion problem. If the model of the 3D scene is known, as in [3], the problem can easily be solved. Otherwise, the problem is very difficult. Theoretically, resolving occlusion could be achieved [9] by inferring a dense map from a stereo pair (or from two consecutive images) so as to compare the depth of the virtual and that of the real object. In fact, despite new advances in 3D reconstruction, the depth map lacks accuracy and cannot be used as is.

Our problem is of course closely related to that of finding occlusions between objects in the real scene. They can be computed by searching disparity discontinuities [4], but this is not an easy problem! Fortunately, we show in this paper that a simpler approach can be used for augmented reality because we have only to discriminate between the known object (i.e. the virtual object) and any object of the scene, and not between unknown objects in the scene.

We then advocate in this paper a contour-based approach that allows us to recover the boundaries of the occlusion mask without 3D reconstruction.

## 2 The link between the occlusion mask and the contours in the images.

Adding a virtual object in the scene requires that the user specifies the relationship of this object with the scene. The



**Figure 1. The occlusion mask for a realistic (a) or an unrealistic (b) position**

user can either give the geometric position of the object in the scene frame, or visually place the object in the image through an interface: the user can then rotate or translate the object in the scene frame until a satisfying position has been reached. In this case, the user does not have at his disposal the 3D map of the scene but only a 2D image (or a sequence), and he can only visually assess the result. Hence, he may sometimes make mistakes while adding an object in the scene. For instance, he may put an object lower than the level ground and unrealistic situations may arise. Consider for instance the problem of adding a virtual mushroom on a real cube behind another cube (Fig.1). If the position of the mushroom is realistic (Fig.1.a) the occlusion mask is only composed of contours present in the scene and of contours belonging to the silhouette of the virtual object (Fig.1.a). For an unrealistic position (see Fig.1.b, the position specified by the user is lower than the top of the cube and the mushroom penetrates the cube), the occlusion mask of the mushroom contains a new contour (in white) that cannot be recovered from the image nor from the silhouette of the virtual object (Fig.1.b). This contour can only be computed if the 3D model of the real scene is known.

Our conclusion is as follows: if the position of the virtual object is realistic, the occlusion mask can be computed using only image data (the edge map) and the silhouette of the virtual objects. Otherwise, the mask contains some contours that can be recovered only if the 3D model of the scene is known.

We think that most of the time, the user can consistently specify the position of the object he wants to add in the scene, possibly with small errors. Indeed, visual assessment generally enables collisions between virtual and real objects to be avoided.

Hence, the occlusion mask is only composed with contours of the image stemming from objects standing in front of the virtual object  $V$  and from parts of the silhouette of the virtual object.

The underlying idea of our algorithm is to identify all

the contours in the image which stand in front of the virtual object without resorting to 3D reconstruction. By grouping the contours according to a proximity criterion, we produce a first rough approximation of the occlusion mask. This first estimation can possibly be somehow erroneous, as all the contours of the occluding objects are not necessarily extracted by the edge detector. Conversely, some edge parts can be mis-identified and can be labeled as *in front of* although they are behind the virtual object. To cope with these problems, we resort to regularization.

### 3 Overview

In this section, we describe the outlines of our algorithm for resolving occlusions using two successive images. Let  $V$  be the virtual object to be added in the real scene. We assume that both the motion between the two frames and the camera parameters are known. In our case, the motion is computed using object based registration.

The main steps of our algorithm are summarized below. For each step, we refer the reader to Fig. 4, which illustrates the key points of the algorithm. In the example shown, we attempt to add a virtual rectangle parallel to the frontal plane of the calibration grid. The motion has been computed using the small circles on the calibration grid.

#### Step 1: Computation of the initial mask

We compute what we call the initial mask  $m_I$ , that is the region in the image to which the object  $V$  corresponds, assuming that the whole object is visible. We thus compute the occluding contours of the virtual object (Fig. 4.b).

#### Step 2: Tracking the contours

The contour chains are extracted in the region of the image corresponding to  $m_I$  (Fig.4.c) (we consider all the chains that cross the initial mask). Then these chains are tracked (Fig.4.d) in the next image using a curved-based tracker which we have developed previously (for further details, see [1])<sup>1</sup>. Finally, the matching of the contours points between the two images is performed by using the epipolar constraint : for each contour point  $m$  belonging to a given chain  $C$ , the corresponding point is computed as the intersection between the epipolar lines and the corresponding tracked chain  $C'$ . As usual, heuristic considerations or constraints are used if several intersections are available.

#### Step 3: Labeling the contour points with *behind* or *in front of*

This important step is detailed in section 4 and is based on the comparison for each image point  $m$  between (i) the 2D displacement associated to the real point in the scene that projects onto  $m$  and (ii) the displacement associated to the 3D point belonging to the virtual object that projects onto

<sup>1</sup> Too small chains are not tracked, as the tracking process is not reliable in this case

$m$  (Fig 4.e and f). The originality of this method is that it does not require 3D reconstruction of the scene.

#### Step 4: Computation of the occlusion mask

This step is detailed in sections 5 and 6. The problem is to recover the occlusion mask from the contours labeled *in front of*. Using the force field created by these contours, the active contours allow us to detect a first estimation of the mask. This estimation is then refined using the initial image (Fig. 4.g,h,i and j).

## 4 Discriminating between “behind” and “in front of” points

### 4.1 The criterion

Fig 2 illustrates the basic geometry of the camera model; for each camera position we consider the coordinate system  $\mathcal{R}_1$  (resp  $\mathcal{R}_2$ ) where the  $z$  axis is the optical axis and the  $(x, y)$  axis are parallel to the image axis.

Let  $I_1$  and  $I_2$  be two consecutive frames. Let  $m_1 \in I_1$  and  $m_2 \in I_2$  be two corresponding points given by the tracking process. We note  $Z_{real}$  the  $z$  coordinate of the real point of the scene corresponding to  $(m_1, m_2)$ . Now let  $M_{obj}$  be the 3D point belonging to the virtual object, whose projection in  $I_1$  is  $m_1$ . This virtual point occludes the real one only if  $Z_{obj} < Z_{real}$ . We show in the following how to solve that without computing  $Z_{real}$ .

Since the calibration is known the projection  $m_{obj}$  of  $M_{obj}$  in  $I_2$  can be computed. Fig. 2 shows intuitively that the points lying in front of or behind the object can be inferred from the relative position of  $m_2$  and  $m_{obj}$  on the epipolar line. More precisely, let us define  $f_{m_1}$  (Fig. 2):

$$f_{m_1} : Z \rightarrow proj_{I_2}(m_{1x}, m_{1y}, Z)$$

where  $m_1 = (m_{1x}, m_{1y})$  and  $proj_{I_2}$  is the projection in image  $I_2$ .

Let  $u_1, v_1, k_{u1}, k_{v1}$  (resp  $u_2, v_2, k_{u2}, k_{v2}$ ) be the intrinsic parameters of the cameras used for the first two frames. The euclidian displacement mapping  $\mathcal{R}_1$  onto  $\mathcal{R}_2$  is called  $[R, T] = \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \end{bmatrix}$ .

Hence, the coordinates in  $\mathcal{R}_1$  of the 3D point  $M$  that projects onto  $m_1$  are

$$M = \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = \begin{pmatrix} \frac{m_{1x} - u_1}{k_{u1}} Z \\ \frac{m_{1y} - v_1}{k_{v1}} Z \\ Z \end{pmatrix} \quad (1)$$

Hence

$$f_{m_1}(Z) = \begin{pmatrix} u_2 + k_{u2} \frac{r_{11}X + r_{12}Y + r_{13}Z + t_x}{r_{31}X + r_{32}Y + r_{33}Z + t_z} \\ v_2 + k_{v2} \frac{r_{21}X + r_{22}Y + r_{23}Z + t_y}{r_{31}X + r_{32}Y + r_{33}Z + t_z} \end{pmatrix} \quad (2)$$

Substituting (1) into (2) allows then  $f_{m_1}(Z)$  to be expressed as an homographic function of  $Z$  whose coefficients depend on the calibration process and on the image point  $m_1$ .

Let us return to the occlusion problem. We have

$$m_2 = f_{m_1}(Z_{real}) \text{ and } m_{obj} = f_{m_1}(Z_{obj})$$

Due to the monotony of the homography, it is therefore easy to compare  $Z_{real}$  and  $Z_{obj}$ : let  $a, b, c, d$  be the coefficients of  $f_m^x$ , that is  $f_m^x(Z) = \frac{aZ+b}{cZ+d}$ . If  $(ad - bc > 0)$ , then

if  $m_2 > \frac{a}{c}$  and  $m_{obj} < \frac{a}{c}$ , then  $Z_{real} < Z_{obj}$

else if  $m_2 < \frac{a}{c}$  and  $m_{obj} > \frac{a}{c}$ , then  $Z_{real} > Z_{obj}$

else if  $m_2 < m_{obj}$  then  $Z_{real} < Z_{obj}$

else if  $m_2 > m_{obj}$  then  $Z_{real} > Z_{obj}$

The two homographies  $f_m^x$  and  $f_m^y$  can be used to decide whether the point is in front of or behind the virtual object. Each point  $m$  can therefore be labeled with *behind* or *in front of* only if the test on the two coordinates give the same result. Otherwise, the label is *doubtful*.

### 4.2 The apparent contours

Contours present in an image can be categorized as either *regular* contours or *apparent* contours [8]. Unlike regular contours that correspond to the same physical event in the 3D scene, apparent contours (or silhouette) are viewpoint dependant (Fig. 3). Hence,  $m_1$  and  $m_2$  are the image of two points  $P_1$  and  $P_2$  lying on the object, but they do not correspond to the same physical point.

The preceding criterion has been established for points belonging to regular contours. Nevertheless, we also use it for points belonging to apparent contours. We explain why in Fig. 3: for an apparent contour, the corresponding points  $m_1$  and  $m_2$  belongs to the tangent lines  $\tau_1$  and  $\tau_2$ . Hence  $m_2$  is the projection of the intersection  $P$  of the two tangent lines. Then, using the above criterion with  $(m_1, m_2)$  leads us to compare the positions of  $P$  and  $V$  instead of comparing the position of  $V$  and  $P_1$ . If  $V$  does not lie between  $P$  and  $P_1$ , the result given by the criterion is correct (outside this interval  $P$  and  $P_1$  are on the same side of  $V$ , i.e. both before or both in front of  $V$ , and the criterion succeeds).

Problems may obviously arise if the virtual object could be added between  $P$  and  $P_1$ . We claim that this is not practically possible for most cases. To prove that, let us compute  $PP_1$ . Let  $R$  be the radius of the real object. We have

$$PP_1 < \widehat{P_1 P_2} = R\alpha$$

where  $\alpha$  is the angle between the tangents; it is usually small as the motion between two frames is small. We can



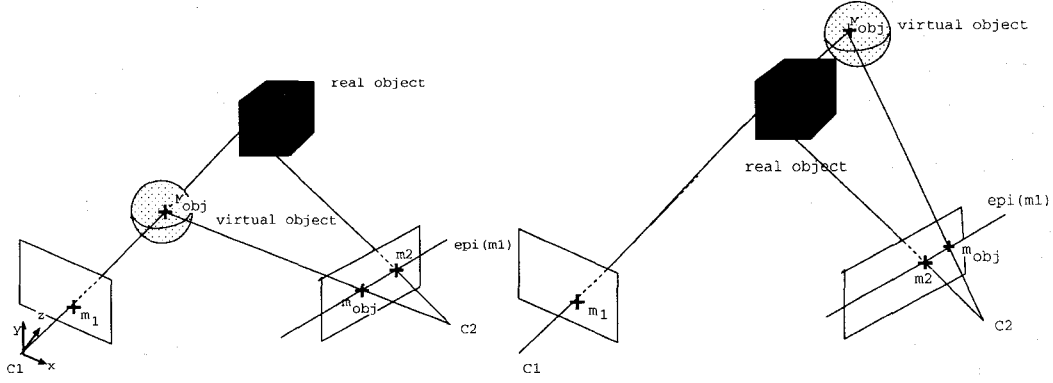


Figure 2. The relative positions of the real and the virtual objects.

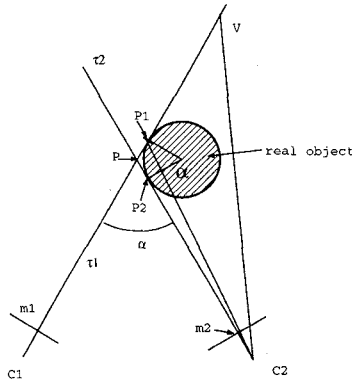


Figure 3. The case of apparent contours.

approximate  $\alpha$  with  $\alpha \approx \arctg(C_1 C_2 / d)$  where  $d$  is the distance between the camera and the object. For our practical case, the distance camera/object is around 2 meters, the distance  $C_1 C_2 < 5cm$  and the radius of the curved parts of the objects in the scene are less than 8cm. Hence  $PP_1 < 8 \times \arctg(5/200) = 0.2cm = 2mm$ . Thus the distance where the criterion fails is small. Moreover, because of visual assessment, the user can hardly attempt to add an object in so small an area.

## 5 Identification of the occluding objects

After the previous step, each contour point is labeled with *behind*, *in front of* or *doubtful*. We now consider the new connected chains  $C_i$  that are only made up of *in front of* points. Our aim is to group the chains that correspond to the same occluding object (Fig. 4.e).

Our algorithm is based on a proximity criterion. As usual, we define the distance of two curves  $C_i$  and  $C_j$  by  $distance(C_i, C_j) = \inf_{x \in C_i, y \in C_j} d(x, y)$ . Given a threshold  $s$ , we consider that two chains such that  $distance(C_i, C_j) < s$  belong to the same object (in prac-

tice  $s$  is equal to a few pixels).

We can therefore build a proximity graph  $\mathcal{G}$ : the nodes represent the chains  $C_i$ ; two nodes are connected only if the distance between the corresponding chains is less than  $s$ .

Then, detecting the occluding objects amounts to compute the cliques in the proximity graph, that is the sets of curves  $\mathcal{H} \subset \mathcal{G}$  such that:  $\forall C \in \mathcal{H}, \exists C' \in \mathcal{H} | C \text{ and } C' \text{ are connected}$ .

## 6 Computing the occlusion masks

Inferring the occlusion mask from a set  $\mathcal{H}$  of contours is an arduous task for several reasons:

- $\mathcal{H}$  may contain some chains that do not belong to the real occlusion mask. This may for instance be due to a threshold  $s$  which has been chosen too large. Problems can also originate in tracking errors.
- We have to build what we call the *envelop* of a set of contours. But all the contours that composed the occluding mask are of course not necessarily detected in the initial contour map and are of course not tracked. Thus, we have to compute the curve that approximates the set  $\mathcal{H}$  as well as possible (Fig 4).

To cope with these problems, we resort to the regularization theory and especially to the active contours models [5]. The underlying idea is to add regularity constraints on the solution that will produce a unique solution, and overcome the problems stemming from lacking or erroneous contours.

It is well known that the snakes naturally tend to shrink if the image gradient is null. We take advantage of this property in our application by initializing a snake with a closed curve outside the set  $\mathcal{H}$  (in practice, we use the smallest rectangle that contains  $\mathcal{H}$  (Fig. 4.g). Then, we let the snake evolve under the influence of the force field created by the contours in the following way. Let  $I_0$  be the contour image:

$$\begin{aligned}
I_0(x, y) &= 255 \text{ if } (x, y) \text{ belongs to one} \\
&\quad \text{of the chains } \in \mathcal{H} \\
&= 0 \text{ otherwise}
\end{aligned}$$

Then let  $F = I_0 * \text{Gauss}(\sigma)$  be the convolution of the contour image with a gaussian kernel with a sufficiently large standard deviation ( $\sigma = 2$  for instance), in order to create the influence field.

As the initializing curve contains  $\mathcal{H}$ , the snake will shrink itself until it reaches the contours of  $\mathcal{H}$ . This will produce the most regular curve resting on  $\mathcal{H}$ . In fact, the curve does not always contain all the chains belonging to  $\mathcal{H}$ . This is quite normal as bumps or spurious contours can be removed by the regularization process.

## 7 Results

Our algorithm has been extensively explained on the example shown in Fig. 4. Another more complex example is shown in Fig. 5, where the virtual object is occluded by two objects of the scene. For the sake of simplicity, the virtual object is a rectangle that is near both the pot and the clown (the distance is around 1 cm, whereas the size of the scene is around 40 cm). Fig. 5.a shows the result of the classification criterion: the points labeled *in front of* are plotted in black, whereas the points labeled *behind* are plotted in gray. The results are visually correct except for the contours that correspond to letters written on the book: these contours are very jagged and are of course difficult to track with sufficient accuracy. This explains the small errors in the classification. Four occluding objects are found after the grouping step (Fig. 5.b). The smallest ones are discarded because the area surrounded by these contours is too small. Using the snake strategies on the other objects allows then the two occluding objects to be detected (Fig. 5.c). The final result is quite convincing.

Another observation can be made from this experiment: all the contour points do not appear on Fig.5.a (see for instance the parallel lines on the pot). This is due to the fact that the contour at these points is almost parallel to the epipolar line (which are approximately horizontal here) and the matching fails. In the neighborhood of such points, the matching is not very accurate; this explains why some parts of the parallel lines of the pot are labeled *behind* instead of *in front of*. Finally, a funny example merging a real scene and strip cartoon heros is shown in Fig.5.d.

## 8 Concluding remarks

We have presented an alternative to the 3D reconstruction for resolving occlusions in augmented reality. The

main interest of our approach is that it only involves 2D computation to check the position of the virtual object against the scene. Resorting to regularization then allows us to compute the mask compatible at best with the classification stage. Moreover, our algorithm is quick and easy to implement.

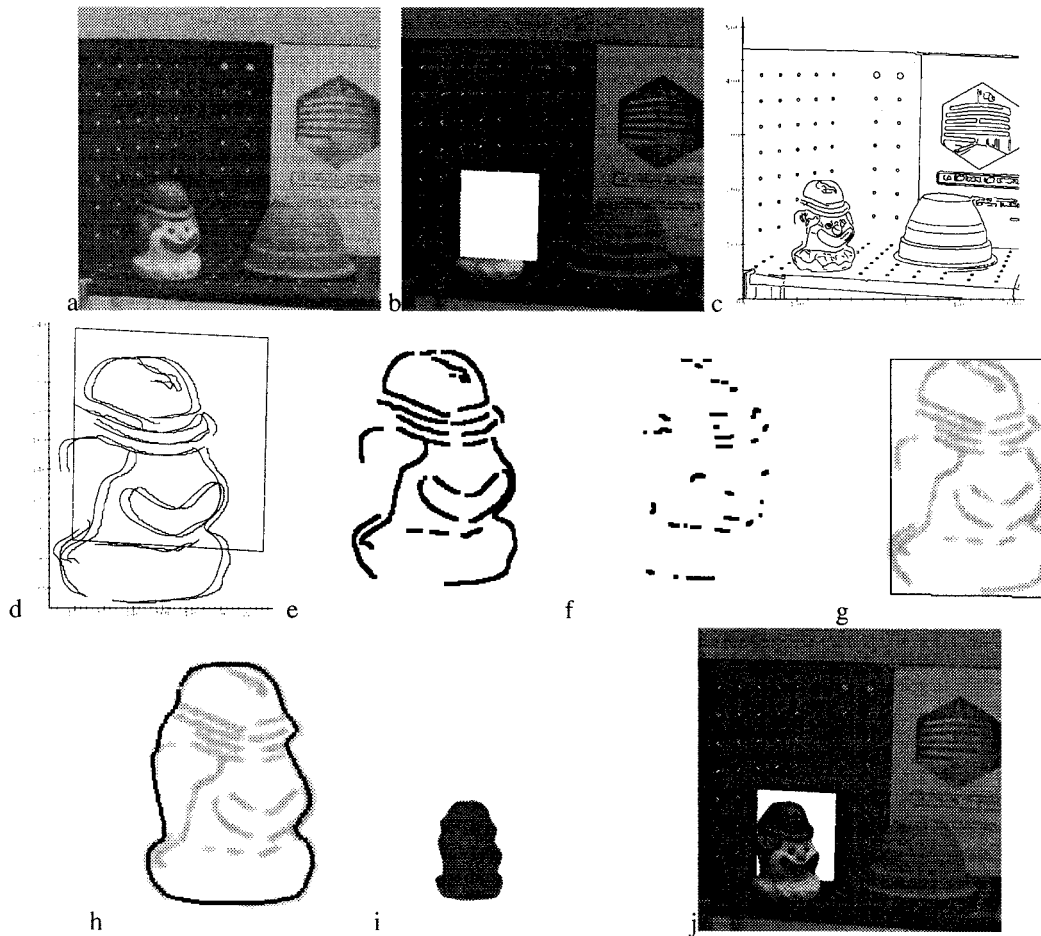
Of course, our approach depends on the quality of the contour map. Some lacking contours may have no influence on the final result as the snake process fills in possible holes in the contours. Conversely, one missing contour may have a great influence on the result. Consider for instance the case of the pot in Fig. 5. As all the interesting contours are not always detected by the edge detector, and as the epipolar matching fails, the pot is almost split into two parts. Fortunately, the upper contour of the pot ensures the connection between the two parts of the pot.

Thus, one missing contour may lead the algorithm to detect two occluding objects instead of one object and may lead us to trouble. Increasing the value of  $s$  is not a satisfactory solution because it will sometimes lead the process to merge possible close occluding objects. Hence, a possible improvement could consist in multi scale analysis, in order to take advantage of some weak contours that do not appear at high scale.

To conclude, even if our method would probably give less interesting results for very textured environments, it is well suited for images where the contour information is relevant.

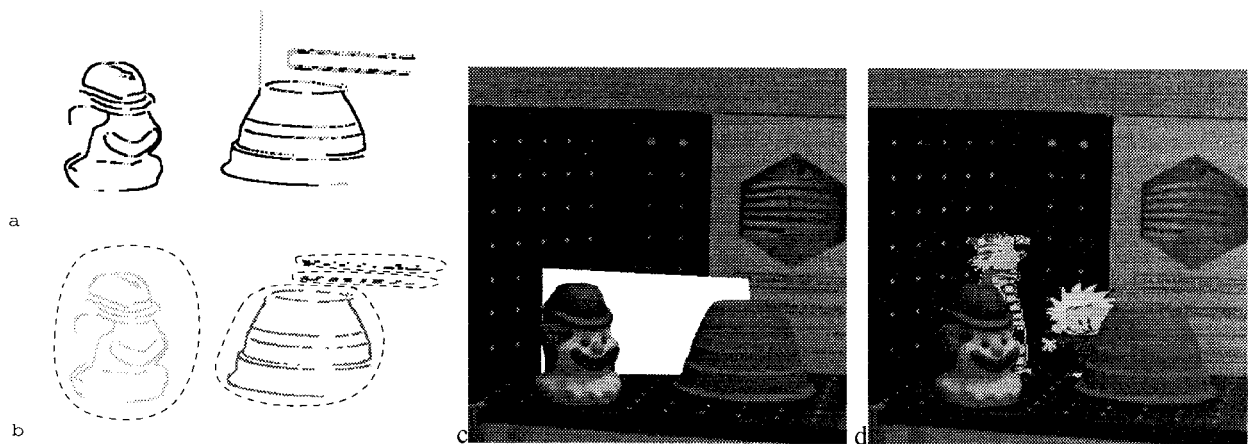
## References

- [1] M.-O. Berger. How to Track Efficiently Piecewise Curved Contours with a View to Reconstructing 3D Objects. In *Proceedings of the 12th ICPR, Jerusalem*, pages 32–36, 1994.
- [2] M.-O. Berger, G. Simon, S. Petitjean, and B. Wrobel-Dautcourt. Mixing Synthesis and Video Images of Outdoor Environments: Application to the Bridges of Paris. In *Proceedings of the 13th ICPR, Vienna*, pages 90–94, 1996.
- [3] D. Breen, R. Whitaker, E. Rose, and M. Tuceryan. Interactive Occlusion and Automatic Object Placement for Augmented Reality. In *EUROGRAPHICS'96, Poitiers, France*, 1996.
- [4] D. Geiger, B. Ladendorf, and A. Yuille. Occlusions and Binocular Stereo. *IJCV*, 14:211–226, 1995.
- [5] M. Kass, A. Witkin, and D. Terzopoulos. Snakes: Active Contour Models. *IJCV*, 1:321–331, 1988.
- [6] S. Ravela, B. Draper, J. Lim, and R. Weiss. Tracking Object Motion Across Aspect Changes for Augmented Reality. In *ARPA Image Understanding Workshop, Palm Spring (USA)*, Aug. 1996.
- [7] M. Uenohara and T. Kanade. Vision based object registration for real time image overlay. *Journal of Computers in Biology and Medicine*, 1996.
- [8] R. Vaillant and O. Faugeras. Using Extremal Boundaries for 3-D Object Modeling. *IEEE Transactions on PAMI*, 14(2):157–173, Feb. 1992.
- [9] M. Wloka and B. Anderson. Resolving Occlusions in Augmented Reality. In *Symposium on Interactive 3D Graphics Proceedings, (New York)*, pages 5–12, Aug. 1995.



**Figure 4. Adding a virtual rectangle in a real scene**

(a) the original image (b) the mask of the virtual object to be added overlaid on the image (c) the edge map (d) result of the tracking process (the two corresponding curves are shown) (e) contours that are labeled *in front of* (f) mis-classified points (g) the gradient field created by the contours labeled *in front of* and the initialing snake (h) the mask after the snake process (i) the mask of the occluding object (j) the result.



**Figure 5. A virtual object occluded by two objects.**

(a) the classification stage: the points in black are *in front of*, the points in gray are *behind* (b) four objects are obtained after the grouping stage (c) the result of the composition (d) a funny composition.