

# An interactive authoring system for video object segmentation and annotation

Huitao Luo<sup>a,\*</sup>, Alexandros Eleftheriadis<sup>b</sup>

<sup>a</sup>*Hewlett-Packard Labs, 1501 Page Mill Road, MS 1203, Palo Alto, CA 94304, USA*

<sup>b</sup>*Department of Electrical Engineering, Columbia University, New York, NY 10027, USA*

Received 4 February 2002; accepted 15 April 2002

---

## Abstract

An interactive authoring system is proposed for semi-automatic video object (VO) segmentation and annotation. This system features a new contour interpolation algorithm, which enables the user to define the contour of a VO on multiple frames while the computer interpolates the missing contours of this object on every frame automatically. Typical active contour (snake) model is adapted and the contour interpolation problem is decomposed into a two-directional contour tracking problem and a merging problem. In addition, new user interaction models are designed for the user to interact with the computer. Experiments indicate that this system offers a good balance between algorithm complexity and user interaction efficiency. © 2002 Elsevier Science B.V. All rights reserved.

*Keywords:* Snake; DP; Graph search; Interactive Rubberband; Authoring tool; MPEG-4

---

## 1. Introduction

To segment and track semantically meaningful video objects (VOs) through a video sequence is currently an interesting topic. The need of VOs arises from two major multimedia applications. One of them is extended video annotation on acquired traditional video sequences. At the editing stage, VOs are specified and annotated with hyperlinks. When the user clicks on the object in any frame that it appears, the annotated data pops up. This is similar to the hyperlink mechanism in HTML. Because video is becoming rapidly

a popular media form, VO annotation is also becoming a very important topic in multimedia authoring research. Actually, several authoring systems are already available that support VO annotation, such as IBM's HotVideo [10] and HyperVideo from Veon [17]. In addition, VO annotation is also an important technique for interactive TV, in which hot links are embedded along with video information [2]. The other application associated with VOs is the new object oriented video representation technique, which is specified within the framework of MPEG-4 and upcoming MPEG-7 standards. In the terminology of MPEG-4, VOs can be segmented from traditional video sequences. Each VO is compressed and decompressed individually. The composition of multiple video objects is coordinated by MPEG-4 systems. Though these two applications are

---

\*Corresponding author. Tel.: +1-650-857-2631; fax: +1-650-857-2951.

*E-mail addresses:* huitao@exch.hpl.hp.com, luoh@ee.columbia.edu (H. Luo), eleft@ee.columbia.edu (A. Eleftheriadis).

different, they have the common problem at the authoring stage, that is, how to specify VOs from large amount of video data in both spatial and temporal domains.

Segmenting objects from image and/or video data has been under intensive research. Numerous algorithms are available in the literature. However, it is becoming widely accepted recently that fully automatic segmentation is difficult. Instead, a semi-automatic approach is a more feasible solution. Several papers have been published according to this idea, for example [3,6,18,9], etc. These papers all sidestepped full automation and designed tracking algorithms to work along with user interaction. Nevertheless, though their contributions to the “tracking” (algorithm) part of the problem are different, their user interaction models are all very simple, i.e., the user defines a VO in the first frame, and then lets the computer track the VO temporally. This model cannot meet the requirements of interactive authoring tools. For example, a common problem is that no quality criteria are proposed for the computer to detect the loss of tracking and thus ask for additional user input. Instead, the user has to observe the tracking course from time to time and offer new input when he or she finds it necessary.

In this paper, we propose a new “interpolation” approach for semi-automatic VO segmentation, and discuss an authoring tool prototype design based on this approach. Our focus is to design the algorithm and the user interaction models at the same time, which helps to solve the efficiency problem in interactive authoring systems. In our system, the user defines a video object by specifying its contour on multiple anchor frames rather than only on the first frame. The computer then uses input information from multiple anchor frames to “interpolate” the VO contours on every frame. Compared with pure tracking approaches, the interpolation approach offers more “interaction points” between the algorithm and the user. From the algorithm’s point of view, the algorithm makes use of user input more efficiently, because each user input contour contributes to VO definitions on those frames before as well as after it, while in the tracking case, a user input only influences the frames after it. From the user’s point

of view, the new approach makes the system performance more predictable because the user can define a VO on frames where large occlusion or motion occurs and most tracking algorithms are likely to fail. In addition, it is more controllable in that with two or more input VO contours, the possible interpolated contours can be effectively limited to certain search regions.

More specifically, the problem can be defined as follows. Given two input contours  $C_b$  and  $C_e$  of a VO on frame  $b$  and frame  $e$ , try to find the object contours  $C_i$  on frames  $i$ ,  $i = b + 1, b + 2, \dots, e - 1$ . We call it a “contour interpolation” problem. As a comparison, we express the “contour tracking” problem as: given an input contour  $C_b$ , try to find the object contour  $C_i$  on frame  $i$ ,  $i = b + 1, b + 2, \dots, e$ . It is natural to consider an interpolation problem as two tracking problems, i.e., to maximize the use of input information on two frames, we can track the input contour from  $C_b$  to  $C_e$  also well as from  $C_e$  to  $C_b$ . In this sense, all the available VO tracking algorithms can be used. However, how to merge the results from these two-directional tracking and produce a final best result is an open problem.

To maximize the use of user input on two frames, we use active contour models (snakes) [11] in our interpolation algorithm. In a snake model, a planar curve is parameterized with nodes and local energy functions defined for each node. The final shape and position of the curve is determined by the global minimization of the snake’s energy. In our work, a traditional snake model is extended in the following ways. First, we use nodes to represent snakes and we design a “contour matching” algorithm to match the node representations of two user input contours. Based on contour matching, contour temporal smoothness and shape similarity criteria are defined. Later discussion will show that these criteria are essential for merging multiple tracking results. Second, we extend the two-dimensional snake model to a three-dimensional model in which new energy terms that reflect spatial temporal constraints are included. Third, a “parametric neighborhood template” is designed to improve the robustness against background edge noises during the active contour tracking course.

The algorithm design in this work is directly influenced by Geiger et al.'s work [8], in which the idea of contour matching was first proposed. In our work, we further develop their contour-matching algorithm by introducing different local motion models. Based on contour-matching results, the interpolation problem is then modeled as a bi-directional snake tracking problem and a merging problem. In addition, similar work on spatial/temporal active snake model can be found in [12,7,4], etc. In Lin and Chang's work [12], efforts were made to combine the active contour model with motion estimation. Quality criteria were also suggested to evaluate the quality of contour tracking. However, their experimental results were not good because no node neighborhood information was used. Chiueh et al. [4] used similar snake technique to track VO contours for annotation. Fu et al. [7] tried to solve the occlusion problem in active contour tracking by segmenting the contours into multiple segments. They used neighborhood information for motion estimation and got better results than [12]. We will later show in this paper that their work can be included in ours as a specific case of our parametric template.

The user interface model in this work is related to the work in [15]. In their system, the user selects key points and the computer grows the corresponding contour segment that links the key points. The current key point is moved by the user until the contour segment grown by the computer is desirable. This process involves both the computer's searching as well as the user's decision. It is an efficient model for user/machine interaction. In our system, an improved active searching mechanism (*Interactive Rubberband*) is designed for the user to define the initial VO contours on anchor frames. In addition, an *iterative interpolation* mechanism is designed for the user to offer error feedback to the interpolation algorithm.

This paper is organized as follows. In Section 2, we introduce the contour representation and contour-matching algorithm. In Section 3, we discuss in detail our contour interpolation algorithm based on contour matching. In Section 4 we summarize the user interaction models employed in the system. Experimental results are presented

in Section 5, and we present concluding remarks in Section 6.

## 2. Contour representation and matching

### 2.1. Contour representation

In this work, a contour can be represented by a vector array  $\{\mathbf{v}_{s,k}\}$  ( $s = 0, 1, \dots, N$ ), where  $k$  is the temporal location and  $s$  is the spatial index of contour pixels. The spatial location of each contour pixel is denoted by vector  $\mathbf{v}_s = (x_s, y_s)$ . In addition, for concise representation and easy matching, a contour can also be represented with subsampled nodes  $\{\mathbf{v}_{S_k(i),k}\}$ , where  $S_k : i \mapsto j$ ,  $i \in [0, 1, \dots, N_s]$ ,  $j \in [0, 1, \dots, N]$  ( $N_s$  and  $N$  are the number of nodes and pixels, respectively) is the subsampling function related to temporal position  $k$ . For example, a uniform subsampling function is defined as  $S_k(i) = i \cdot \text{unit}$ , where unit is the uniform subsampling length. In this paper, we name  $\{\mathbf{v}_{s,k}\}$  *pixel representation* and  $\{\mathbf{v}_{S_k(i),k}\}$  *node representation*. Note that the node representation is actually a first-order polynomial approximation of the pixel representation.

### 2.2. Contour matching

The purpose of contour matching is to find the correspondence between two contours  $C_b$ ,  $C_e$ , which are the contours of the same VO at the different temporal locations ( $b$  and  $e$ ). This is essential for interpolating the object contours between them. In the pixel representation, the length of  $C_b$  and  $C_e$  are not necessarily the same and a pixel-by-pixel correspondence cannot be created. Therefore, we use subsampled node representation for contour interpolation study, and a contour-matching algorithm is employed to find the correspondence between the two subsampled contours.

Mathematically, the matching process can be defined as follows. Given two input contours:  $C_b = \{\mathbf{v}_{s,b}\}$ ,  $C_e = \{\mathbf{v}_{s,e}\}$  in the pixel representation, and a subsampled node representation for the first contour  $\{\mathbf{v}_{S_b(i),b}\}$ , find the corresponding node representation for the second contour  $\{\mathbf{v}_{S_e(i),e}\}$ .

Here the matching of the nodes of two contours can be expressed with the mapping function  $f_m : \mathbf{v}_{S_b(i),b} \mapsto \mathbf{v}_{S_e(i),e}$ ,  $i \in [0, 1, \dots, N_s]$ . Generally, we fix the node representation of the first contour  $C_b$  by uniform subsampling (other subsample algorithms are also possible, see [16] for a detailed discussion), and the matching process is reduced to finding the corresponding subsample function  $S_e(i)$  for the second contour.

In order to find the mapping function  $f_m$ , we define a local energy term for each matched node pair. The final matching result is determined by the global minimization of the matching energy of the two contours. This is similar to the matching approach in [8]. In this work, we assume the motion of the considered VO is nonrigid globally, but rigid locally, and the shape of its two contours is similar locally everywhere. This assumption is true in general if the motion of the VO is not too fast in relation to the temporal distance between the two frames on which the contours  $C_b$  and  $C_e$  are defined. A number of rigid motion models, i.e., translation and affine, are possible for local matching energy definition.

*Translation motion model.* If we denote the motion vector between two matched nodes as  $MV_i = \mathbf{v}_{S_b(i),b} - \mathbf{v}_{S_e(i),e}$ , then the matching energy term is defined as

$$E_i = \mu \|MV_i - MV_{i-1}\| + \eta (S_e(i) - S_e(i-1))^2,$$

where the first term is the smoothness evaluation of the motion vectors of two neighboring nodes, the second term is an elastic constraint on the distance of two neighboring nodes, and  $\mu, \eta$  are two weighting factors.

*Affine motion model.* In contrast to the translation model, the affine model needs three motion vectors to define. Here we denote the affine mapping function as  $\mathcal{A}_i = \text{Affine}_i(MV_{i-1}, MV_i, MV_{i+1})$ . Under this function, the local contour segment  $\{\mathbf{v}_{s,b}\}$ ,  $s \in [S_b(i-1), S_b(i+1))$  is mapped to a pixel set  $\{\mathbf{v}'_e(s)\}$  on frame  $e$ , while its corresponding contour segment on frame  $e$  is denoted as pixel set  $\{\mathbf{v}_{S_e(s),e}\}$ ,  $s \in [S_e(i-1), S_e(i+1))$ . Then the matching energy term is defined as

$$E_i = \mu D(\{\mathbf{v}'_e(s)\}, \{\mathbf{v}_{S_e(s),e}\}) + \eta (S_e(i) - S_e(i-1))^2,$$

where the operator  $D(\cdot)$  is the distance measure of two sets, which we define as

$$D(A, B) = \left[ \sum_{\mathbf{v}_i \in B} \min_{\mathbf{v}_j \in A} \|\mathbf{v}_i - \mathbf{v}_j\| + \sum_{\mathbf{v}_j \in A} \min_{\mathbf{v}_i \in B} \|\mathbf{v}_j - \mathbf{v}_i\| \right] / 2. \quad (1)$$

In practice, (1) can be implemented in an iterative manner for each pixel as follows. We define  $d(\mathbf{v}_i, A) = \min_{\mathbf{v}_j \in A} \|\mathbf{v}_i - \mathbf{v}_j\|$ , and denote  $d_n(\mathbf{v}_i, A)$  as the value of  $d(\mathbf{v}_i, A)$  in the  $n$ th iteration. At the initialization stage

$$d_0(\mathbf{v}_i, A) = \begin{cases} 0 & \text{if } \mathbf{v}_i \in A, \\ +\infty & \text{otherwise.} \end{cases}$$

The iteration is then defined as

$$d_{n+1}(\mathbf{v}_i, A) = \min_{\mathbf{v}_j \in N(\mathbf{v}_i)} (d_n(\mathbf{v}_j, A) + \|\mathbf{v}_i - \mathbf{v}_j\|),$$

where  $N(\mathbf{v}_i)$  is the 8-neighborhood set of pixel  $\mathbf{v}_i$ .

With the definition of local matching energy terms, the matching problem is converted to an energy minimization problem. This can be easily solved by DP algorithm [1], which we do not discuss in detail here. In practice, the affine model is more complex than the translation model, but the quality is better, especially when the subsample distance between neighboring nodes is big. In Fig. 1, we show a result of contour matching under the translation model. Fig. 1(a) is the 50th, Fig. 1(b) is the 70th frame of the Carphone sequence. The green lines are the contours and the red lines link the matched node pairs.



Fig. 1. Results of contour matching for the Carphone sequence: (a) is the 50th and (b) is the 70th frame of the Carphone sequence. The green lines are the user specified contours and the red lines link the matched node pairs.

### 3. Contour interpolation algorithm

#### 3.1. Localized energy minimization model

The essence of Kass et al.'s snake model [11] is to define a local energy term for each contour node, and the shape of the contour is determined by minimizing the total snake energy globally. When we go from a two-dimensional spatial snake to three-dimensional spatial/temporal snake, it is possible to extend the local energy terms from 2D to 3D as well. In this work, we study the extended local energy terms as intraframe energy terms and interframe energy terms, respectively. From now on, because no subsample issue will be involved, snakes are assumed to be in the node representation. The node representation notation  $\{\mathbf{v}_{S_k(i),k}\}$  is simplified to  $\{\mathbf{v}_{i,k}\}$  whenever possible.

*Intraframe energy.* As usual, the intraframe energy terms include two parts, a gradient term and a smoothness term,

$$E_{\text{intra},i} = \eta_{\text{edge}} E_{\text{gradient},i} + \eta_{\text{smooth}} E_{\text{smooth},i}, \quad (2)$$

where the first gradient term is defined as

$$E_{\text{gradient},i} = \int_{s=S(i-1)}^{S(i)} \frac{255}{(10 + \|\nabla(\mathbf{c}(\mathbf{v}_s))\|)} ds,$$

in which 255 and 10 are two empirical values.  $\mathbf{c}(\mathbf{v})$  is the color vector of node  $\mathbf{v}$ , and the second smoothness term is defined as

$$E_{\text{smooth},i} = \left( \frac{2\|\mathbf{v}_{S(i-1)} + \mathbf{v}_{S(i+1)} - 2\mathbf{v}_{S(i)}\|}{\|\mathbf{v}_{S(i-1)} - \mathbf{v}_{S(i)}\| + \|\mathbf{v}_{S(i)} - \mathbf{v}_{S(i+1)}\|} \right)^2,$$

which is designed to eliminate the influence of the node distance on the contour smoothness measure.

*Interframe energy.* In available papers [12,7] on temporal active contour tracking, interframe energy terms generally employed include optical flow, motion smoothness, interframe color, etc. A basic problem in these approaches is that most of their node energy definitions are based only on image features at the node's position rather than on its neighborhood. This makes the color and especially motion information generally not accurate. Ideally, the contour nodes' neighborhood should be observed in order to track them from frame to frame. However, a problem here is, different from typical point tracking problems,

contour nodes are most likely located on the boundary of a moving object, so their neighborhood is not constant. In order to capture the consistency through the tracking course, we introduce a concept of *parametric neighborhood template*. A parametric template  $T$  is defined as a data structure including two arrays: a vector array  $\{\mathbf{dv}_0, \mathbf{dv}_1, \dots, \mathbf{dv}_n\}$  and a weighting array  $\{w_0, w_1, \dots, w_n\}$ , in which a vector  $\mathbf{dv}_i$  represents the position of the  $i$ th pixel in the neighborhood and  $w_i$  represents its contribution weight to template measurements. For each contour node  $\mathbf{v}_{i,k}$ , a parametric template  $T_{i,k}$  is defined and kept updated frame by frame through the tracking course.

With the definition of parametric template  $T_{i,k}$ , the color of two temporally neighboring contour nodes  $\mathbf{v}_{i,k}, \mathbf{v}_{i,k-1}$  can be compared as

$$\begin{aligned} \text{Diff}_c(\mathbf{v}_{i,k}, \mathbf{v}_{i,k-1}, T_{i,k}) \\ = \sum_{\mathbf{dv}_j \in T_{i,k}} w_j \|\mathbf{c}(\mathbf{v}_{i,k} + \mathbf{dv}_j) - \mathbf{c}(\mathbf{v}_{i,k-1} + \mathbf{dv}_j)\|, \end{aligned} \quad (3)$$

where  $\mathbf{c}(\mathbf{v})$  is the color vector of node  $\mathbf{v}$ . In addition, we define the motion vector at node  $\mathbf{v}_{s,k}$  as

$$\begin{aligned} \text{MV}(\mathbf{v}_{i,k}, T_{i,k}) \\ = \sum_{\mathbf{dv}_j \in T_{i,k}} w_j \text{MV}^{(p)}(\mathbf{v}_{i,k} + \mathbf{dv}_j) / \sum_{\mathbf{dv}_j \in T_{i,k}} w_j, \end{aligned} \quad (4)$$

where  $\text{MV}^{(p)}(\mathbf{v}_{i,k})$  is the estimated motion vector at  $\mathbf{v}_{i,k}$ . If we do not consider occlusion, a simple way to determine the weights of template  $T$  is to set  $w_j$  for those neighboring pixels inside the contour as 1 and for those outside the contour as 0. This can be illustrated in Fig. 2. For the occlusion case as was discussed in [7], we can easily switch the weights by setting inside weights to 0 and outside weights to 1. In general cases, both the size and the weights of the parametric template can be adjusted flexibly to get the best results of the contour node tracking.

Based on parametric template, the interframe energy terms for each contour node are enumerated as follows:

1. Color similarity:  $E_{\text{color},i,k} = \text{Diff}_c(\mathbf{v}_{i,k}, \mathbf{v}_{i-1,k}, T_{i,k})$ , where function  $\text{Diff}_c(\cdot)$  is defined in (3).

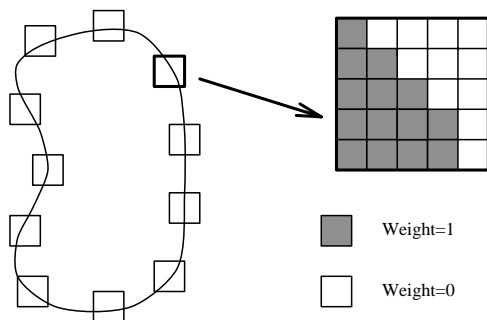


Fig. 2. Illustration of parametric template concept. During the tracking course, both the size and weights of the template can be adapted.

2. Optical flow:  $E_{\text{optical},i,k} = \text{MV}(\mathbf{v}_{i,k-1}, T_{i,k-1}) - \text{MV}_{i,k-1}$ , where the first term  $\text{MV}(\mathbf{v}_{i,k-1}, T_{i,k-1})$  is defined in (4) and we assume forward motion estimation is used. The second term is  $\text{MV}_{i,k-1} = (\mathbf{v}_i - \mathbf{v}_{i,k-1})$ .
3. Motion smoothness:  $E_{\text{motion},i,k} = \|\text{MV}_{i-1,k} + \text{MV}_{i+1,k} - 2\text{MV}_{i,k}\|$ . This is a smoothness measure for motion vectors on the spatially neighboring nodes. It is minimized if the local motion is in the translation form.
4. Shape stiffness:  $E_{\text{shape},i,k} = \|\text{angle}(\mathbf{v}_{i-1,k-1}, \mathbf{v}_{i,k-1}, \mathbf{v}_{i+1,k-1}) - \text{angle}(\mathbf{v}_{i-1,k}, \mathbf{v}_{i,k}, \mathbf{v}_{i+1,k})\|$ , where  $\text{angle}(\mathbf{v}_{i-1}, \mathbf{v}_i, \mathbf{v}_{i+1})$  is the angle based on the three spatially neighboring nodes. This term measures the local shape similarity. It is minimized if the local motion is rigid.
5. Temporal smoothness:  $E_{\text{temporal},i,k} = \|\mathbf{v}_{i,k-1} + \mathbf{v}_{i,k+1} - 2\mathbf{v}_{i,k}\|$ . This is the smoothness measure for the three temporally neighboring nodes.

The interframe energy  $E_{\text{inter}}$  is then the weighted sum of above terms.

*Search algorithm for minimization.* With the definition of local energy terms, the contour interpolation problem can be expressed as an energy minimization problem as follows. Given two contours  $\{\mathbf{v}_{i,b}\}, \{\mathbf{v}_{i,e}\}$ , ( $b < e$ ,  $i = 0, 1, \dots, N_s$ ), find the contour nodes  $\{\mathbf{v}_{s,k}\}$  ( $k = b + 1, \dots, e - 1$ ,  $s = 0, 1, \dots, N_s$ ) that minimize the global energy  $\sum_{s,k} (E_{\text{inter},s,k} + E_{\text{intra},s,k})$ .

Though it is natural to extend the local energy terms from 2D to 3D, the increase in computational complexity is an important problem. In the

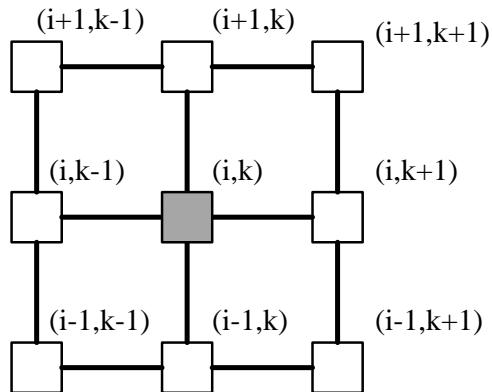


Fig. 3. Spatial temporal neighborhood of a contour node for the local energy definition. In the figure,  $i$  is the spatial index and  $k$  is the temporal index.

2D case, each node  $\mathbf{v}_i$ 's local energy is defined in relation to two neighbors, i.e.,  $E_{\text{intra}} = f(\mathbf{v}_{i-1}, \mathbf{v}_i, \mathbf{v}_{i+1})$ , while in the 3D case, the local energy terms for each node  $\mathbf{v}_{s,k}$  is defined in relation to eight neighbors! This change is illustrated in Fig. 3. If each node has a search region of size  $n$ , the local search complexity then increases from  $n^3$  to  $n^9$ , which makes global minimization algorithm difficult to design. Though powerful algorithms such as *simulated annealing* should still be able to solve this minimization problem, the slow speed of convergence makes it inappropriate for an interactive segmentation tool.

In this work, a sub-optimal solution is obtained by converting the contour interpolation problem into two tracking problems: a forward tracking from  $C_b$  to  $C_e$  and a backward tracking from  $C_e$  to  $C_b$ . They are also referred to as bi-directional tracking in this paper. When converted to a tracking problem, the neighborhood definition of the current pixel  $(i, k)$  cannot include nodes in a neighboring frame that has not been processed. Therefore, we shift the neighborhood definition in Fig. 3 by one frame. For example in the forward tracking model, if the current node is  $(i, k)$  in Fig. 3, then its eight nodes are  $(i + 1, p), (i, p), (i - 1, p)$ ,  $p = k - 2, k - 1, k$  (not including  $(i, k)$ ). Among them, six are already fixed and only two variable nodes  $(i + 1, k)$  and  $(i - 1, k)$  will influence its local energy. This is the same as the case with

the intraframe energy  $E_{\text{intra},s}$ . Therefore, it is easy to design a search algorithm with DP. After the bi-directional tracking, another search process is used to find the optimal contours out of the previous tracking results.

### 3.2. Bi-directional tracking

In the tracking model, each node  $\mathbf{v}_{i,k}$ 's total node energy can be written as

$$E_{\text{total},i,k} = f_E(\mathbf{v}_{i-1,k}, \mathbf{v}_{i,k}, \mathbf{v}_{i+1,k}) \quad (5)$$

because the other six neighboring nodes  $\mathbf{v}_{i+p,k-2}$ ,  $\mathbf{v}_{i+p,k-1}$ ,  $p = (-1, 0, 1)$ , in the previous frames are all fixed. This energy expression is similar to those used for 2D active contour models, except that the detailed expression of  $f_E$  is different. Therefore, we use the DP algorithm similar to that used in [1].

*Limited search region versus searching complexity.* According to [1], the local energy expression in (5) is a second order expression in the sense that it includes two neighboring nodes as variables. In this case, a two-element vector  $(\mathbf{v}_{i+1}, \mathbf{v}_i)$  is used as the status index for DP. The DP search is then carried out as follows:

$$S_i(\mathbf{v}_{i+1}, \mathbf{v}_i) = \min_{\mathbf{v}_{i-1}} [S_{i-1}(\mathbf{v}_i, \mathbf{v}_{i-1}) + f_E(\mathbf{v}_{i-1}, \mathbf{v}_i, \mathbf{v}_{i+1})].$$

Note that the frame index  $k$  is omitted in the above expression because no temporal information is involved.<sup>1</sup> If the search size for each node is  $n$  and the total number of nodes in each contour is  $m$ , the complexity of DP is  $O(mn^3)$ . Obviously, the search size  $n$  is an important factor in the overall complexity and should be limited as much as possible.

Two clues are used to limit the search size in our work. First at the global level, a search stripe can be created for each matched node pair. This is illustrated in Fig. 4. For ease of discussion, the temporal orbit of the matched node is projected into one frame. On this frame, a search stripe is defined. Note that the definition of a search stripe is different based on different spatial location of the matched node pairs. In Fig. 4, two basic type of search stripe definitions are depicted. The

orientation of width and height are defined differently as well based on the different orientations of the matched node pairs. In practice, both the width and height of the global search stripe can be determined by the user in an interactive way, according to the motion of VO. That is, if the motion is more like a pure translation, the height of the stripe  $S_{\text{height}}^{(G)}$  may be reduced, otherwise it is increased. The bottom line is that the global search region should be at least big enough to cover the temporal orbit of the node's motion. In addition, in the case of large search stripes that exceed a certain threshold, the global search stripe is further re-sampled spatially to reduce the overall search complexity. The two axes "x" and "y" used for resampling are marked in Fig. 4. Note that the parallelogram definition of search region is more efficient than a rectangular one in the computational sense, while without much performance degradation.

Once the global stripes are defined for node pairs, the tracking of nodes is constrained within their stripes on every frame. In addition, at the local level, the local search region is further determined frame by frame by the forward motion vector at the current frame. In our work, the determination of  $S_{\text{height}}^{(L)}$  and  $S_{\text{width}}^{(L)}$  is based on  $S_{\text{height}}^{(G)}$  and  $S_{\text{width}}^{(G)}$ , and the local search is carried out on top of the re-sampled grids created for the global search stripe, i.e. the re-sampled grids along the  $x$ - and  $y$ -axis.

In above two clues, the global stripe limits the possible location of the local search region. Note that in the pure interpolation case, both the width and height of the global search stripe is zero. Therefore, the global search stripe is actually a generalization of interpolation.

*Closed contour problem.* Until now, the discussed contours  $C_i$  are all by default open. In practice when contours are constrained to be closed, the minimization search used for tracking purpose can be approximated with two-pass open contour searching. That is, for a closed contour  $C_i$ ,  $i = 0, 1, 2, \dots, n$ , and  $C_0 = C_n$ , its minimization status can be found as follows. First break the closed contour at node  $C_0$ , use previous algorithm to process open contour  $C_i$ ,  $i = 0, 1, 2, \dots, n-1$ , which produces a temporary contour  $C'_i$ . Second,

<sup>1</sup>Note the temporal information is implicitly used in the energy term  $f_E(\mathbf{v}_{i-1}, \mathbf{v}_i, \mathbf{v}_{i+1})$ .

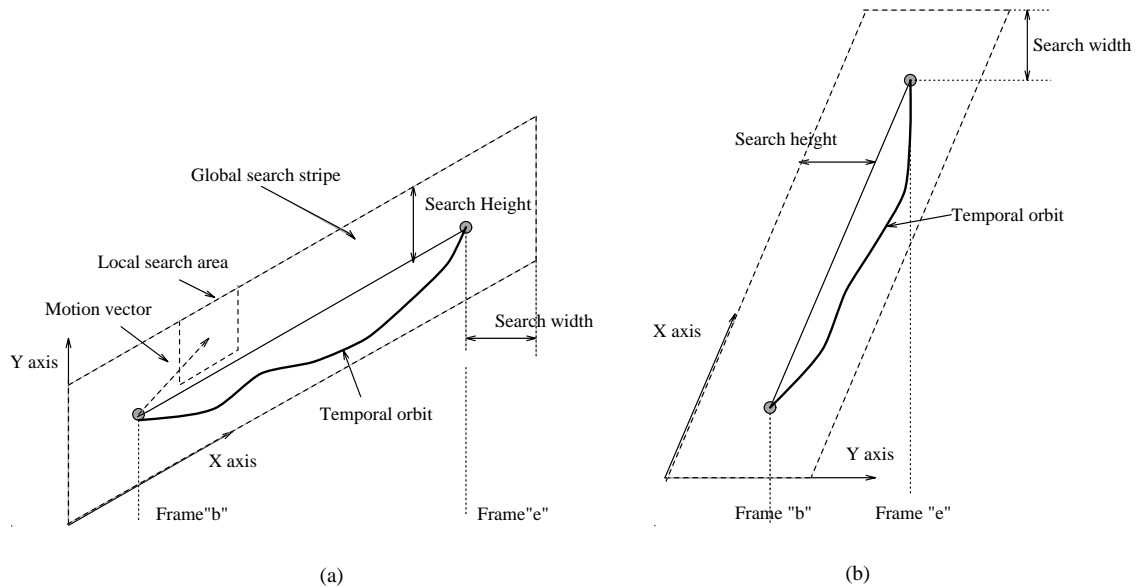


Fig. 4. Global and local search region limitation. The two grayed circles represent the matched node pair. The size of the global search region is determined by “search height” and “search width”, while the local search region is determined by motion estimation.

close  $C'_i$  by setting  $C'_n = C'_0$  and then break it half way at node  $C'_{(n/2)}$ , this produces an open contour  $C''_i$ .  $C''_i$  is further processed with the open contour algorithm and the final result is obtained.

### 3.3. Merging of multiple results

Though the bi-directional tracking approach reduces the search complexity, its limitation is that it only makes use of user input in one frame (either  $C_b$  or  $C_e$ ) at a time. Due to the error accumulation, the tracked contour  $C_k$  always degrades when  $k$  approaches  $e$ , if tracked from  $C_b$  to  $C_e$ , and vice versa. Actually we have observed that if the object’s motion involves self-occlusion and/or uncovering, sometimes it is very difficult for the active contour model to track its contour in one direction, but easy to do it in the other direction. Fig. 5 shows such an example. In Fig. 5(a) and (b) are the user defined VO on the 118th and the 132nd frame of the Carphone sequence. From frame 118 to frame 132, the man’s left ear is uncovered because of the rotation of his head. If the object contour is tracked forward, i.e. from the 118th frame to the 132nd frame, the uncovered ear was not included as part of the VO. This is

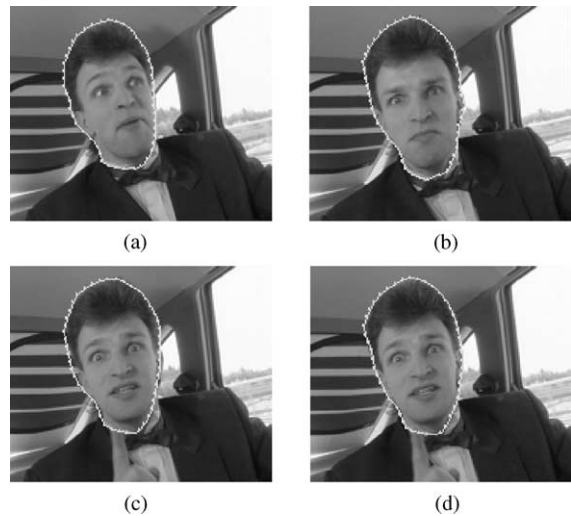


Fig. 5. Illustration of the necessity of bi-directional tracking and result merging: (a) and (b) are the user defined VO on the 118th and the 132nd frame of the Carphone sequence; (c) is the tracked object contour on the 125th frame by forward tracking; (d) is the tracked object contour on the 125th frame by backward tracking.

depicted in Fig. 5(c). On the other hand, if the contour is tracked backward, i.e. from the 132nd frame to the 118th frame, the motion is reversed



and the uncovering motion of left ear changes to occlusion, which is easy for the active contour algorithm to handle. The result of backward tracking for frame 125 is shown in Fig. 5(d). Note that the left ear is correctly included as part of the video object.

Therefore, a good algorithm to merge the results from the two tracking processes is important. In this work, an efficient DP algorithm is designed for the merging purpose. The problem can be defined as follows. Given two set of contours  $\{C_k^{(1)}\}, \{C_k^{(2)}\}$  ( $k = b, b + 1, \dots, e$ ) that are created by two contour tracking processes (one from  $C_b$  to  $C_e$  and the other from  $C_e$  to  $C_b$ ), find a contour set  $C_k$  ( $k = b, b + 1, \dots, e$ ,  $C_k = C_k^{(1)}$  or  $C_k = C_k^{(2)}$ ) that meets certain merit criteria as the final output of contour interpolation. In the terminology of DP, we can say that the target is to find an optimal path from  $C_b$  to  $C_e$  that maximize the merit criteria.

In this work, the merit criteria for each candidate contour include two terms: a temporal smoothness energy term  $E_T$ , and a shape merit energy term  $E_S$ , i.e.,

$$E_{\text{merge}}(C_k) = \eta_T E_T(C_k) + \eta_S E_S(C_k). \quad (6)$$

The first term is defined in a localized form,

$$E_T(C_k) = \sum_{i=0}^{i=N_s} \|\mathbf{v}_{i,k-1} + \mathbf{v}_{i,k+1} - 2\mathbf{v}_{i,k}\|, \quad (7)$$

where  $N_s$  is the number of nodes in each contour. Note that  $E_T(C_k)$  is different from previously defined  $E_{\text{temporal}}$  in that  $E_T(C_k)$  is defined for each

contour while  $E_{\text{temporal}}$  is defined for each contour node.

The second term of (6) is based on the shape similarity comparison between two contour pairs:  $(C_b, C_k)$  and  $(C_k, C_e)$ . If we denote the shape similarity measure of two contours  $C_k$  and  $C_l$  as  $\text{shape}(C_k, C_l)$ , then the  $E_S$  term can be written as

$$E_S(C_k) = [w_1(k) \text{shape}(C_b, C_k) + w_2(k) \text{shape}(C_e, C_k)], \quad (8)$$

where  $w_1(\cdot)$  and  $w_2(\cdot)$  are two weighting functions. They are designed as a linear function of frame indices  $k, b$  and  $e$ . In addition, the shape similarity measure used here can be defined based on two interframe node energy terms  $E_{\text{shape}}$  and  $E_{\text{motion}}$ , as discussed in Section 3.1. At the contour level, this expression is defined as

$$\text{shape}(C_k, C_l) = \sum_{i=0}^{i=N_s} \eta_1 [(\mathbf{v}_{i,k} - \mathbf{v}_{i,l}) - (\mathbf{v}_{i-1,k} - \mathbf{v}_{i-1,l})] + \eta_2 [\text{angle}(\mathbf{v}_{i-1,k}, \mathbf{v}_{i,k}, \mathbf{v}_{i+1,k}) - \text{angle}(\mathbf{v}_{i-1,l}, \mathbf{v}_{i,l}, \mathbf{v}_{i+1,l})].$$

In above (6)–(8),  $E_{\text{merge}}(C_k)$  is actually defined in relation to three contours:  $C_{k-1}$ ,  $C_k$  and  $C_{k+1}$ . Therefore, it is easy to solve the minimization problem with DP. Here Fig. 6 is used to illustrate the DP based merging algorithm. In Fig. 6, each circle represents a contour  $C_k^{(d)}$ , DP is used to find an optimal path in the temporal direction that has the best merit quality, i.e., in the sense of temporal smoothness and shape similarity.

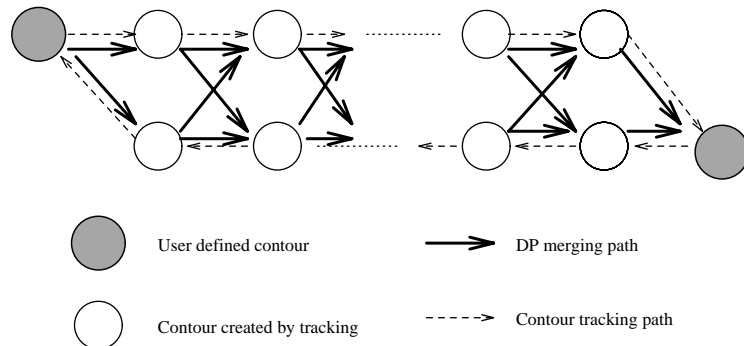


Fig. 6. Illustration of DP approach for merging the bi-directional contour tracking. Each circle represents a contour  $C_k^{(d)}$ , DP is used to find an optimal path in the temporal direction that has the best merit quality.

#### 4. User interaction model

In a typical semi-automatic system, the user's role includes two important functions. One is to give the initial data for the computer to begin the computation, the other is to correct the computer's errors. In our system, these two functions are handled by *Interactive Rubberband* and *Iterative Interpolation*, respectively.

##### 4.1. Interactive Rubberband

In available image authoring systems, two types of solutions are common for a user to specify an object contour in an image (or a video frame). In one of them (Type One), the user specifies every point of the contour, while the computer does nothing but record the positions of each mouse click and links the positions with line segments. Typical examples include the polyline drawing in XFIG and free style drawing in PHOTOSHOP. This type of solution gives the user full control of the shape and position of the contour, but ignores the computational power of the computer. Obviously, it is tedious to input an accurate contour point by point. On the other hand, in the other type of solutions (Type Two), the image is modeled as grids and the contour as paths linking the grids. The user has only to select a starting point and an ending point of a contour segment, the computer finds the whole segment by searching the minimal cost path that links the two points. This can be done with either dynamic programming or graph search algorithms such as Dijkstra's algorithm [5]. Publications belonging to this type include [13–15], etc. Compared with Type One solutions, Type Two approach relieves the user's labor by introducing computer searching during the interaction. However, its problem is that the user has less control on the contour. Sometime when the gradient information within the image is complex, an intended contour segment may be attracted to an erroneous strong neighboring edge, which is totally undesirable. In addition, the "Active-Scissors" approach defined in [15] requires the computer to calculate the optimal path to every pixel within the image every time a new contour point position is chosen by the user. This

is not efficient if the size of image is big and real time performance is hard to achieve.

To overcome the problems while retaining the benefits of above two groups of solutions, we design an Interactive Rubberband tool, which is a hybrid of the two of them.

An Interactive Rubberband is in principle a graph search based edge detection algorithm that is similar to the aforementioned Type Two solution. The difference is that it comes with an adjustable containing rectangle that limits the range of graph search. This is illustrated in Fig. 7. The user moves the current moving point, which, together with the fixed point, determines a containing rectangle. Graph search is then carried out within the rectangle and a contour segment is grown to link these two points. Compared with the Active-Scissors approach in [15], Interactive Rubberband is more efficient because it limits the graph search range to the neighborhood of the desirable contour segment. Moreover, the user can control the result of graph search by adjusting the width of the rectangle, e.g. if there is strong noisy edges in the neighborhood, the user may get rid of them by narrowing the containing rectangle. In the extreme case, the width can be set to one, then the Interactive Rubberband reduces to above type one solution. In this sense, the Interactive Rubberband is a generalization of the previous Types One and Two solutions.

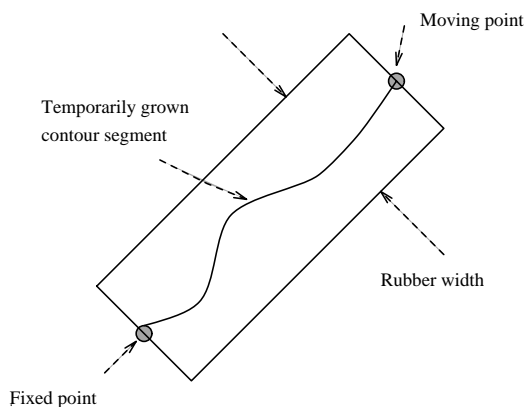


Fig. 7. Illustration of active rubberband. The containing rectangle is determined by two points: a fixed point and a moving point, and the rubberband width. The width of the rubberband is adjustable by the user.

### 4.2. Iterative Interpolation

Though in experiments the discussed interpolation algorithm showed good performance, error is inevitable in practice, especially when the two anchor frames on which the user specifies object contours  $C_b$  and  $C_e$  are far away in the temporal direction. Iterative interpolation is found to be a good solution to this problem.

In general, the reasons for errors in interpolation are complex. However, in the interpolation algorithm based on bi-directional tracking, a video object contour can always be reliably tracked from  $C_b$  to a certain  $C_{(b+t_1)}$ , and from  $C_e$  to a  $C_{(e-t_2)}$ , where  $t_1 > 0$  and  $t_2 > 0$ . If  $b + t_1$  turns out equal to  $e - t_2$ , the problem is solved. Otherwise, we can move the  $C_b$  to  $C_{(b+t_1)}$ , and  $C_e$  to  $C_{(e-t_2)}$ , and begin the interpolation again. In this way, the interpolation is done iteratively until the two contours  $C_b$  and  $C_e$  converge. This process can be better illustrated in Fig. 8. In Fig. 8, points  $P(1)$  and  $P(2)$  are a matched node pair on contours  $C_b$  and  $C_e$ . After the first round of interpolation, point  $P(1)$  is correctly tracked to  $P(3)$  and  $P(2)$  to  $P(4)$ . At this stage, the user changes the  $C_b$  to the temporal position at  $P(3)$  and  $C_e$  to the position at  $P(4)$ , sets the global search parameters accordingly, and begins the interpolation again. As depicted in the figure, the global search area in

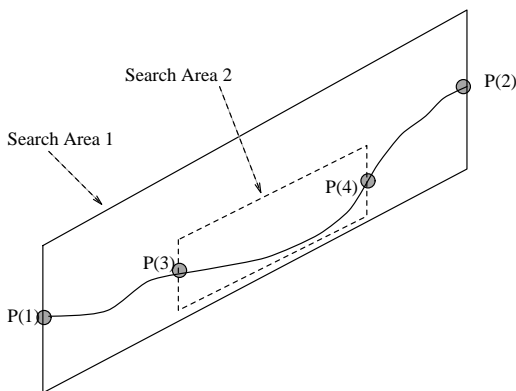


Fig. 8. Illustration of iterative interpolation. Point  $P(1)$  and  $P(2)$  are a matched node pair on initial contours  $C_b$  and  $C_e$ . After the first round of interpolation, point  $P(1)$  is correctly tracked to  $P(3)$  and  $P(2)$  to  $P(4)$ . Their corresponding contours are used as new  $C_b$  and  $C_e$  and the interpolation is done iteratively until converges.

the second round for point pair  $P(3)$ – $P(4)$  is much smaller than that of  $P(1)$ – $P(2)$ . This is an important factor that helps the iterative interpolation process converge.

In Fig. 9, a practical example is given on the Foreman sequence to show how the iterative

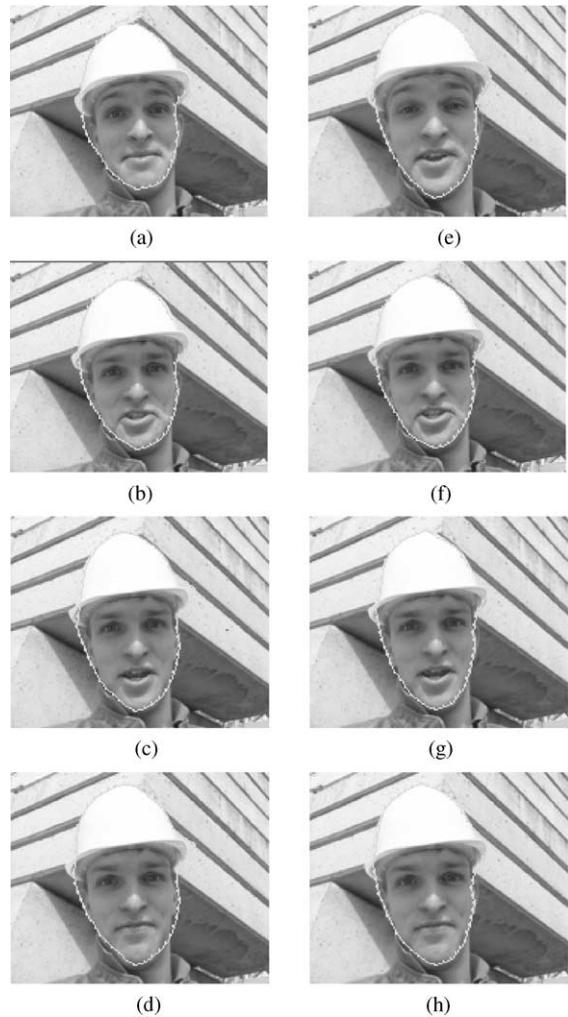


Fig. 9. An example of iterative interpolation on the Foreman sequence. In the first round, the user specifies  $C_b$  at frame 50 and  $C_e$  at frame 100, the global search parameters are height = 15(pixels), width = 4(pixels). (a)–(e) are results on frames 80, 84, 87, 90, 94 at this round. In the second round, previous contour result on frame 80 is used as  $C_b$  and contour result on frame 94 is used as  $C_e$ , the global search parameters are reduced to be height = 4, width = 4. The new interpolation results on frame 84, 87, 90 are showed in (f)–(h), respectively.

interpolation works. In the first interpolation round, the user specifies  $C_b$  on frame 50 and  $C_e$  on frame 100, the global searching parameters are



Fig. 10. Comparison of the effect of parametric template on active contour tracking. The left column is the tracking result without and the right column is the result with parametric template. The first row is the beginning user input contour in frame 0, the following two rows are the tracked results for the 43rd and 47th frames.

height=15(pixels), width=4(pixels). Figs. 9(a)–(e) are results on frames 80, 84, 87, 90, 94 in this round. Obviously, the tracking result is poor from frame 80–94, mainly because a strong neighboring edge has attracted the snake erroneously (due to space limits, other frames are not included in the figure). To overcome the error, the user moves  $C_b$  to frame 80 and  $C_e$  to frame 94, and change the global searching parameters to height = 4, width = 4. Based on the contours on frame 80 and 94, the results after the new interpolation round on frames 84, 87, 90 are shown in Figs. 9(f)–(h), respectively. Obviously, the second interpolation round improves the accuracy of tracked contours if we compare the contours in Figs. 9(b)–(d) with those in Figs. 9(f)–(h). It is worth noticing that in the second round of interpolation, the user does not have to tell the computer laboriously what exactly a correct contour is. Instead, the user just chooses new anchor frames on which the tracked contours are correct, and changes the global searching parameters accordingly (limits the global searching area as much as possible). It is the computer's work to do the interpolation again based on new user input information.

## 5. Experiments

The discussed VO annotation system is implemented on PCs running Windows 95. Experiments



Fig. 11. Illustration of tracking results merging. The user defined VO contours on frame 75 and 180. The first row is the result of backward tracking from frame 180 to frame 75, and the second row is the interpolation results that merged bi-directional trackings. The frame numbers, from left to right, are 170, 160, 150, 130 and 90.

are carried out over MPEG-4 testing video sequences as well as sequences from a library used by Columbia's VideoQ<sup>2</sup> system.

First, we compare the effect of parametric template for active contour tracking on Carphone sequence. In the experiment, single-direction forward tracking is used. In Fig. 10, the left column is the tracking result without and the right column is the result with the parametric template. The first row is the beginning user input contour in frame 0, the following two rows are the tracked results for the 43rd and 47th frames. Obviously, the parametric template improves the tracking robustness in complex boundary conditions.

Fig. 11 shows the results of tracking result merging on the Mother–Daughter sequence. The user defines VO contours on frame 75 and 180. The first row is the result of backward tracking from frame 180 to frame 75, and the second row is the interpolation results that merge bi-directional trackings. The frame numbers, from left to right, are 170, 160, 150, 130 and 90. We can see that from frame 180 to 160, the merged results are taken from backward tracking, while from 160 to 75, the merged results are taken from forward tracking (which is not shown here due to the space limit). That is, the merged interpolation results are better than tracking results on either direction. To achieve (approximately) the same segmentation quality, traditional single-directional tracking needs one or two additional user specified initialization of object segmentation on the same testing video.

Fig. 12 is a fully finished segmentation result on the first 100 frames of the Foreman sequence. Figs. 12(a)–(h) are the produced contours on frame 10, 20, 30, 40, 60, 70, 80, 90. To finish the segmentation, the user specified three initial contours on frame 0, 50 and 100. An additional iteration is involved in the interpolation between frame 0 and frame 50, and frame 50 and frame 100, respectively.

In practice, the computational complexity of the algorithm depends on the size of the searching area and the complexity of the contours. In our experiment, a 200-MHz Pentium is used, the

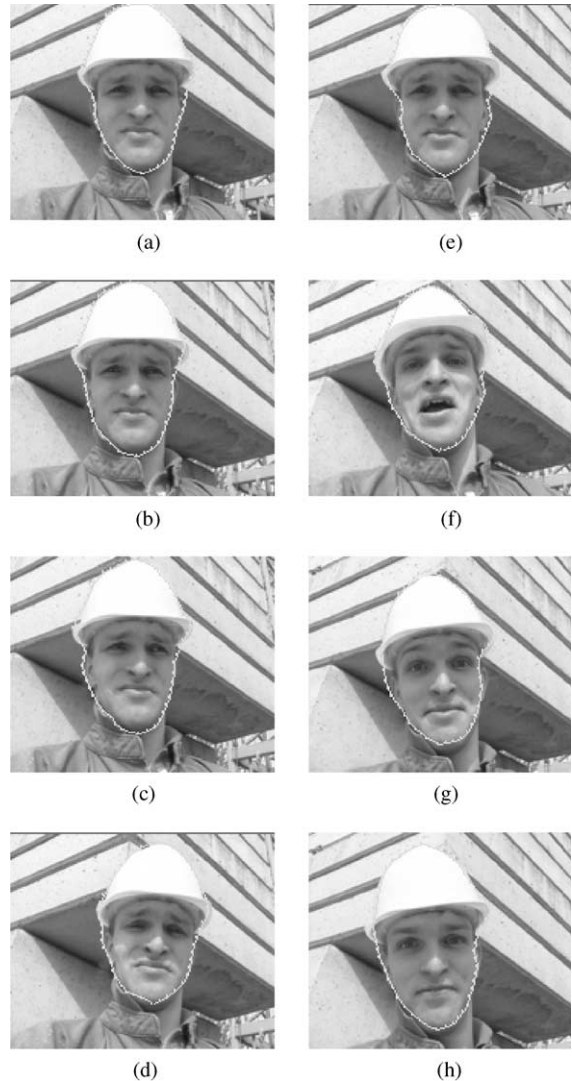


Fig. 12. Illustration of a full segmentation result on the Foreman sequence. The frame numbers are 10, 20, 30, 40, 60, 70, 80, 90 for (a)–(h).

average speed of the interpolation algorithm is about 0.01 s/node and/or 0.6 s/frame (tested on several MPEG-4 sequences in QCIF size).

## 6. Concluding remarks

In this paper, an interactive authoring system is designed for VO segmentation and annotation.

<sup>2</sup><http://www.ctr.columbia.edu/VideoQ>.

This system features a new contour interpolation algorithm, which makes better use of user inputs and has more stable performance than traditional single-directional tracking algorithms. In addition, efficient user interaction models are built for both initial data input and machine error feedback. Our experiments show that this prototype system works efficiently both from the machine's and the user's point of view, in that it balances the user's decision-making capability with the computer's processing and searching power. For the future work, we are looking to explore better ways to model the object motion based on user specification on multiple frames and improve the segmentation performance.

## References

- [1] A.A. Amini, T.E. Weymouth, R.C. Jain, Using dynamic programming for solving variational problems in vision, *IEEE Trans. Pattern Anal. Machine Intell.* 12 (9) (1990) 867–885.
- [2] V.M. Bove Jr., J. Dakss, E. Chalom, S. Agamanolis, Hyperlinked television research at MIT media laboratory, *IBM System J.* 39 (3–4) (2000) 470–478.
- [3] E. Chalom, V.M. Bove Jr., Segmentation of an image sequence using multi-dimensional image attributes, in: *Proceedings of the IEEE International Conference on Image Processing*, Lausanne, September 1996.
- [4] T.C. Chiueh, T. Mitra, C.K. Yang, Zodiac: a history-based interactive video authoring system, in: *ACM International Multimedia Conference*, Bristol, England, September 1998.
- [5] T.H. Cormen, C.E. Leiserson, R.L. Rivest, *Introduction to Algorithms*, MIT Press, Cambridge, MA, 1990, Chapter 25.2.
- [6] P. Correia, F. Pereira, The role of analysis in content-based video coding and indexing, *Signal Processing* 66 (1998) 125–142.
- [7] Y. Fu, A.T. Erdem, A.M. Tekalp, Occlusion adaptive motion snake, *IEEE Trans. Image Process.* 9 (12) (December 2000) 2051–2060.
- [8] D. Geiger, A. Gupta, L.A. Costa, J. Vlontzos, Dynamic programming for detecting, tracking and matching deformable contours, *IEEE Trans. Pattern Anal. Machine Intell.* 17 (3) (1995) 294–302.
- [9] C. Gu, M.C. Lee, Semantic video object tracking using region-based classification, in: *Proceedings of the IEEE International Conference Image Processing*, Chicago, October 1998.
- [10] IBM HotVideo website, <http://www.software.ibm.com/net.media/hotvideo/index.html>.
- [11] M. Kass, A. Witkin, D. Terzopoulos, Snakes: Active contour models, *Int. J. Comput. Vision* 1 (4) (1988) 321–331.
- [12] Y.T. Lin, Y.L. Chang, Tracking deformable objects with the active contour model, in: *IEEE International Conference On Multimedia Computing and Systems*, Ottawa, Canada, June 1997.
- [13] A. Martelli, A heuristic search methods to edge and contour detection, *Commun. ACM* 19 (2) (1976) 73–83.
- [14] E. Mortensen, W. Barrett, Intelligent scissors for image composition, in: *Proceedings of ACM SIGGRAPH 95*, Los Angeles, CA, 1995, pp. 191–198.
- [15] E.N. Mortensen, W.A. Barrett, Interactive segmentation with intelligent scissors, *Graphical Models Image Process.* 60 (1998) 349–384.
- [16] K. Sobottka, I. Pitas, Segmentation and tracking of faces in color images, in: *Proceedings of Second International Conference on Automatic Face and Gesture Recognition*, Killington, Vermont, 1996, pp. 236–241.
- [17] Veon website, <http://www.veon.com/>.
- [18] D. Zhong, S.F. Chang, AMOS: an active system for MPEG-4 video object segmentation, in: *Proceedings of the IEEE International Conference on Image Processing*, Chicago, October 1998.