



following procession. It is implemented using CFD with two preliminary steps, namely, global motion estimation and compensation, change detection.

### 1) Global Motion Estimation

To estimate the global motion, the camera motion is modeled by the six-parameter affine motion model:

$$\begin{cases} x' = ax + by + e \\ y' = cx + dy + f \end{cases} \quad (1)$$

the global motion is estimated iteratively using Gauss-Newton (GN) algorithm. The initial parameters of global motion estimation are obtained by least square method. To reduce the computation time, the GN algorithm is applied within a three-level multi-resolution pyramid, which is generated using [1/4, 1/2, 1/4] filter. For global motion compensation the bilinear interpolation is used.

### 2) Change Detection

Let  $I_t$  and  $I_{t'}$  denote two consecutive frames,  $d_{t,t'}$  is the frame difference:

$$d_{t,t'}(p) = W \times I_t(p) - W \times I_{t'}(p) \quad (2)$$

The change detection mask  $D_{t,t'}$  is defined as:

$$D_{t,t'}(p) = \begin{cases} 1 & \text{if } d_{t,t'}(p) > T \\ 0 & \text{else} \end{cases} \quad (3)$$

where  $W$  is a smooth filter (e.g.  $3 \times 3$  Gaussian filter),  $T$  is the threshold depending on the camera noise, which can be selected within [5, 10] according to the specific occasion of video application.

Connected components analysis [9] is applied to eliminate those smaller noise regions in  $D_{t,t'}$ .

### 3) Continuous Frame Difference (CFD)

There are always some occluded or disoccluded background areas left in  $D_{t,t'}$ . To obtain more accurate initial areas, the CFD is calculated. Let  $D_{t,t-1}$ ,  $D_{t,t+1}$  denote the change detection masks of current frame  $I_t$  and previous frame  $I_{t-1}$ , next frame  $I_{t+1}$ , the CFD mask  $D_t$  is defined as:

$$D_t = D_{t,t-1} \cap D_{t,t+1} \quad (4)$$

By calculating the CFD, more accurate initial areas  $IF_t$  are separated from current background  $IB_t$ , and there are few background areas left in the  $IF_t$ .

## 3. Spatial Segmentation

Watershed algorithm based on the gradient image [7] is applied to segment the initial areas  $IF_t$  due to its robustness and affectivity. To improve the accuracy of segmentation, the  $IF_t$ 's gradient in the  $YCbCr$  color space is calculated.

Let  $G'_Y$ ,  $G'_{Cb}$  and  $G'_{Cr}$  denote the normalized gradient images of the three color components ( $Y$ ,  $Cb$  and  $Cr$ ), which are calculated using Canny's gradient approximation [9], and then normalized into [0, 255]. The color gradient  $G_{col}$  is calculated as:

$$G_{col}(p) = \begin{cases} \max\{\omega_Y \cdot G'_Y(p), \omega_{Cb} \cdot G'_{Cb}(p), \omega_{Cr} \cdot G'_{Cr}(p)\} & \text{if } D_t(p) = 1 \\ 0 & \text{else} \end{cases} \quad (5)$$

where  $\omega_Y, \omega_{Cb}, \omega_{Cr}$  are the weight coefficients of the gradient magnitudes of  $Y, Cb$  and  $Cr$ . In experiments, we have used  $\omega_Y = 0.5, \omega_{Cb} = \omega_{Cr} = 0.25$ .

Watershed is implemented using fast immersion simulation [7] on the initial areas  $IF_t$ . To overcome the watershed's inherent drawback, namely, over segmentation, a spatio-temporal merging scheme [5] is adopted to merge small regions.

For the convenience of narration in the rest of the paper, let  $R_t = \{R_t^1, \dots, R_t^K\}$  denote the set of regions

after spatial segmentation.  $N_i$  is the size of  $R_t^i$ ,  $Nor(R_t^i)$  is the set of neighbors of  $R_t^i$ .  $IF_t = \bigcup_{i=1}^K R_t^i$  is the

initial areas obtained by CFD,  $IB_t = I_t - IF_t$  is the current background.  $IO_{t-1} = \bigcup_{L(R_{t-1}^i) \in F} R_{t-1}^i$  is the objects

segmented from previous frame, where  $L(R_{t-1}^i)$  is the classification of  $R_{t-1}^i$  ( $L(R_{t-1}^i) \in \{F, B\}$ ,  $F$  denotes foreground,  $B$  denotes background).

## 4. Region Classification

MRF model is a very prominent stochastic model applied comprehensively in image processing and computer vision. For region classification, MRF is defined using spatial, temporal and adjacent constraints based on the similarity between the region and current background, the region and the segmented results of previous frame, the region and its neighbors. Then, regions are classified by estimating the MAP of the MRF.

According to the Hammersley-Clifford theorem [8] and Bayes' rule, the MAP of MRF can be estimated by minimizing the posterior energy [5], [8]. So we define the posterior energy  $U_p(X | O)$  of MRF

as:

$$U_p(X | O) = \sum_{i=1}^K \alpha \cdot V_i^S(X, O) + \beta \cdot V_i^T(X, O) + \sum_{(i,j) \in E} \gamma \cdot V_{ij}^R(X, O) \quad (6)$$

where  $V_i^S(X, O)$ ,  $V_i^T(X, O)$ ,  $V_{ij}^R(X, O)$  are the energy functions corresponding to spatial, temporal

and adjacent constraints,  $\alpha, \beta, \gamma$  are the associated weight coefficients,  $E = \{(i, j) | R_i^j \in \text{Nor}(R_i^i)\}$  is the set of all adjacent relationship.

1) spatial constraint energy  $V_i^s(X, O)$

$$V_i^s(X, O) = \begin{cases} f(SD(R_i^i), T_s, SD_{\max}, SD_{\min}) & X_i = B \\ 1 - f(SD(R_i^i), T_s, SD_{\max}, SD_{\min}) & X_i = F \end{cases} \quad (7)$$

where  $SD(R_i^i)$  is the spatial similarity of  $R_i^i$  and  $IB_i$ :

$$SD(R_i^i) = \min_v \frac{1}{N_b} \sum_{l=1}^3 \omega_l \cdot \sum_{p \in R_i^i} |I_i^l(p) - I_i^l(p+v)|, \quad (8)$$

$$D_i(p+v) = 0, N_b > \frac{2}{3} N_i$$

$I_i^1(p), I_i^2(p), I_i^3(p)$  are the associated intensity functions of  $Y, Cb, Cr$ ,  $\omega_1, \omega_2, \omega_3$  are associated weight coefficients same as  $\omega_r, \omega_{cb}, \omega_{cr}$  in section 3.  $v$  is the matching vector within a  $w \times w$  window.  $N_b$  is the number of background pixels matching with  $R_i^i$ .  $f(d, T, d_h, d_t)$  is a scaling function which normalizes  $d$  to  $[0, 1]$ :

$$f(d, T, d_h, d_t) = \begin{cases} 0.5 \times (d - d_t) / (T - d_t) & \text{if } d < T \\ 0.5 + 0.5 \times (d - T) / (d_h - T) & \text{else} \end{cases} \quad (9)$$

$V_i^s(X, O)$  represents the likelihood of the region  $R_i^i$  to be classified as foreground/background based on the magnitude of the spatial similarity  $SD(R_i^i)$  and threshold  $T_s$ . In  $IF$ , only a little part of regions are non-object regions and very similar to nearby areas in  $IB$ , thus, the likelihood of region to be classified as foreground/background can be estimated according to the spatial similarity. Doing this way, the algorithm can overcome the drawback of region classification based on motion estimation which is sensitive to irregular movement and illumination.

2) temporal constraint energy  $V_i^t(X, O)$

$$V_i^t(X, O) = \begin{cases} f(TD(R_i^i), T_t, TD_{\max}, TD_{\min}) & X_i = F \\ 1 - f(TD(R_i^i), T_t, TD_{\max}, TD_{\min}) & X_i = B \end{cases} \quad (10)$$

where  $TD(R_i^i)$  is the temporal similarity of  $R_i^i$  and  $IO_{i-1}$ :

$$TD(R_i^i) = \min_{R_{i-1}^j} \sum_{l=1}^3 \omega_l \cdot | \text{avg}_{p \in R_i^i} I_i^l(p) - \text{avg}_{p \in R_{i-1}^j} I_{i-1}^l(p) |, \quad (11)$$

$$L(R_{i-1}^j) = F$$

$V_i^t(X, O)$  represents the likelihood of the region  $R_i^i$  to be classified as foreground/background based

on the magnitude of temporal similarity  $TD(R_i^i)$  and threshold  $T_t$ . The more similar  $R_i^i$  and  $IO_{i-1}$  are, the more likely  $R_i^i$  is to be classified as foreground. By considering the temporal constraint, mistaken classifications of regions similar to the background can be prevented.

3) adjacent constraint energy  $V_{ij}^R(X, O)$

$$V_{ij}^R(X, O) = \begin{cases} (RD(R_i^i, R_j^j) - RD_{\min}) / (RD_{\max} - RD_{\min}) & X_i = X_j \\ 1 - (RD(R_i^i, R_j^j) - RD_{\min}) / (RD_{\max} - RD_{\min}) & X_i \neq X_j \end{cases} \quad (12)$$

where  $RD(R_i^i, R_j^j)$  is the adjacent similarity of  $R_i^i$  and its neighbor  $R_j^j$ :

$$RD(R_i^i, R_j^j) = \sum_{l=1}^3 \omega_l \cdot | \text{avg}_{p \in R_i^i} I_i^l(p) - \text{avg}_{p \in R_j^j} I_j^l(p) |, \quad (13)$$

$$R_j^j \in \text{Nor}(R_i^i)$$

$V_{ij}^R(X, O)$  represents that the more similar  $R_i^i$  and its neighbor  $R_j^j$  are, the more likely  $R_i^i$  and  $R_j^j$  belong to the same class.

The constants  $\alpha, \beta, \gamma$  determine the relative proportion of the three terms in the posterior energy. In experiments, we found that using  $\alpha = 1.0$ ,  $\beta = \gamma = 0.6$  can obtain satisfactory results. The thresholds  $T_s$  and  $T_t$  depend on the complexity of video's background. To choose accurate thresholds adaptively, we select the spatial similarity and temporal similarity as classification plane respectively, and choose the plane that maximize the distance between classes as the value of  $T_s$  and  $T_t$  according to the Fisher linear discrimination criterion. The minimization of posterior energy is performed using an iterative deterministic relaxation scheme known as HCF [8].

## 5. Experimental Results

Experiments have been carried out on a Pentium IV 2.4G PC, the frame size is  $352 \times 288$ , the processing speed per frame is about 1.5~2.0 s, the average processing time for each step is shown in Table I. If the objects segmentation is performed on the whole image, only the time spent on the classification will exceed 10 s. This shows that temporal segmentation before spatial segmentation and classification reduces the computation complexity largely.

Table I Average Processing Time for Each Step

Steps	Temporal segment	Spatial segment	Classify	Total
Time (ms)	450	482	798	1730

Segmentation results of a diving sequence are shown in Fig. 1. In this sequence, the illumination is obvious and the motion of the athlete is rapid and non-rigid, experiments results demonstrate that the proposed algorithm is not sensitive to objects' irregular movement and illumination.

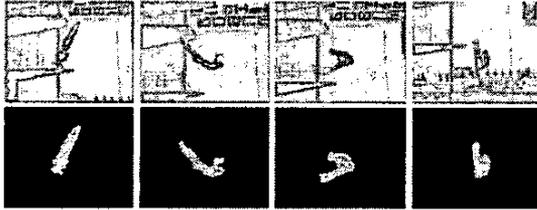


Fig. 1 Segmentation results of Diving sequence

The contrastive experimental result on the 53th, 82th, 125th and 142nd frame of Stepan sequence are giving in Fig. 2. (a) is the results of the algorithm proposed by Y. Tsai [5] using region classification based on motion estimation, (b) is the results of our algorithm. According to the contrastive experiments, we can see that the region classification based on motion estimation results in wrong classification of some disoccluded background regions, however, the algorithm proposed in this paper can overcome the drawback of the region classification based on motion estimation.

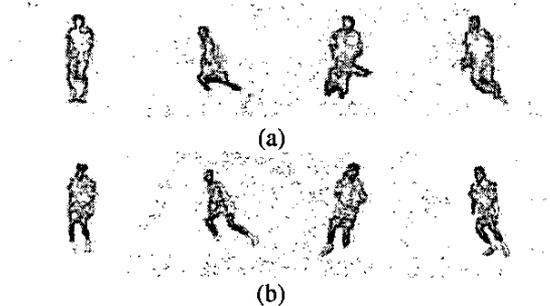


Fig. 2 contrastive experiments

## 6. Conclusions

An automatic segmentation algorithm for moving objects in video sequences under multi-constraints is proposed. Compared with existing hybrid segmentation algorithms, the proposed algorithm has the following advantages. First, temporal segmentation before spatial segmentation and classification reduces the computation time of

segmentation process largely. Secondly, the region classification is implemented by estimating the MAP of the MRF using spatial, temporal and adjacent constraints. Doing this way, the proposed algorithm overcomes the drawback of region classification based on motion estimation, which is sensitive to objects' irregular movement and illumination. Experimental results validate the efficiency of the proposed algorithm.

## 7. References

- [1] D. Wang, Unsupervised video segmentation based on watersheds and temporal tracking, *IEEE Trans. Circuits and Systems for Video Technology*, 1998, 8(5): 539-546.
- [2] Jianping fan *et al*, Automatic image segmentation by integrating color-edge extraction and seeded region growing, *IEEE Trans. Image Processing*, 2001, 10(10): 1454-1466.
- [3] Til Aach, Andre Kaup, "Statistical model-based change detection in moving video," *Signal Processing*, vol. 31, pp. 165-180, 1993.
- [4] A. Neri, S. Colonnese, G. Russo and P. Talone, "Automatic moving object and background separation," *Signal Processing*, vol. 66, pp. 219-232, 1998.
- [5] Y. Tsai and A. Averbuch, "Automatic Segmentation of Moving Objects in Video Sequences: A Region Labeling Approach," *IEEE Trans. Circuits and Systems for Video Technology*, 2002, 12(7): 597-612.
- [6] I. Patras, E. A. Hendriks and R. L. Lagendijk, "Video Segmentation by MAP Labeling of Watershed Segments," *IEEE Trans. Pattern Analysis and Machine Intelligence*, pp. 326-332, 23(3), March 2001.
- [7] L. Vincent and P. Soille, Watersheds in digital spaces: an efficient algorithm based on immersion simulations, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 1991, 13: 583-598.
- [8] P. B. Chou and C. M. Brown, "The theory and practice of Bayesian image labeling," *Int. J. Comput. Vis.*, vol. 4, pp. 185-210, 1990.
- [9] M. Sonka, V. Hlavac and R. Boyle, *Image Processing, Analysis, and Machine Vision*, International Thomson Publishing, 1998.