

PAPER

Real-Time Tracking of Multiple Moving Object Contours in a Moving Camera Image Sequence

Shoichi ARAKI^{†*}, Regular Member, Takashi MATSUOKA^{†**}, Nonmember,
Naokazu YOKOYA[†], and Haruo TAKEMURA[†], Regular Members

SUMMARY This paper describes a new method for detection and tracking of moving objects from a moving camera image sequence using robust estimation and active contour models. We assume that the apparent background motion between two consecutive image frames can be approximated by affine transformation. In order to register the static background, we estimate affine transformation parameters using LMedS (Least Median of Squares) method which is a kind of robust estimator. Split-and-merge contour models are employed for tracking multiple moving objects. Image energy of contour models is defined based on the image which is obtained by subtracting the previous frame transformed with estimated affine parameters from the current frame. We have implemented the method on an image processing system which consists of DSP boards for real-time tracking of moving objects from a moving camera image sequence.

key words: *moving camera image, motion estimation, robust estimation, active contour model, realtime object tracking*

1. Introduction

Motion tracking is one of the most important issues in computer vision. Especially, it is highly difficult to detect moving objects from image sequences obtained by a moving camera since there exists the apparent motion of a static background. A number of methods of motion detection from a moving camera have been proposed; some are based on direct detection of camera parameters with sensors [1], optical flow estimation [2] and registering the static background using a geometric transformation [3]–[5]. Geometric transformation based methods are promising for real-time processing because of their low computational costs. In most existing methods based on registering the static background, it is assumed that the apparent motion of the static background is dominant in an image sequence, and that the motion between consecutive two image frames can be represented by a geometric transformation such as affine transformation. In these methods, all pixels are selected for estimating affine transformation parameters. Therefore, the estimated parameters are not so accurate since the pixels selected for estima-

tion are from not only the background but also moving object regions. It is necessary for estimating accurate parameters to select pixels only from the background. Pixels selected from moving objects are outliers for the parameter estimation. Sawhney et al. have estimated apparent background motion based on the differences of pixel intensity between consecutive two image frames using M-estimation which is a kind of robust estimator [8]. This method is sometimes unstable for estimating motion parameters when an initial solution is not properly given. Moreover, the method is not suitable for real-time processing because of high computational costs caused by recursive calculation in M-estimation. In the field of image coding, parametric motion estimation based on parameterized block correlation [9] has been proposed, but it also needs high computational costs caused by calculating partial differential coefficient at many times.

In this paper, we propose a new method for real-time detection and tracking of moving objects from a moving camera image sequence using robust statistics [6] and active contour models [7]. In the method, we assume that the apparent motion of the static background is dominant in an image sequence and that the background motion between two consecutive image frames can be represented by affine transformation. For registering the static background with small computational time, we estimate affine parameters using LMedS (Least Median of Squares) method based on the location of some feature points and their correspondences [10], [11]. LMedS is a kind of robust estimator useful for a data set including outliers and takes less computational cost compared with M-estimation. Moving object regions can be detected by subtracting the previous frame transformed with estimated affine parameters from the current frame. It is well-known that active contour models are promising for tracking moving objects. We have recently developed split-and-merge contour models for tracking multiple moving objects [12], [13]. In this paper, the image energy term of split-and-merge contour models is defined based on the subtracted image so that contour models are attracted only to edges of moving objects. We have implemented the method on an image processing system which consists of DSP boards for real-time tracking. We have successfully detected and tracked walking persons in

Manuscript received May 27, 1999.

Manuscript revised December 15, 1999.

[†]The authors are with Nara Institute of Science and Technology, Ikoma-shi, 630-0101 Japan.

*Presently, with Matsushita Electric Industrial Co., LTD., Kyoto-fu, 619-0237 Japan.

**Presently, with Mitsubishi Electric Corporation, Fukuyama-shi, 720-8647 Japan.

outdoor scenes by processing each image frame in 1/7.5 second.

This paper is structured as follows. In Sect. 2, we describe the proposed method. In Sect. 3, we explain the implementation of the method on a DSP-based image processing system and show experimental results of real-time tracking. Section 4 summarizes the present work.

2. Global Motion Compensation and Tracking of Moving Objects by Split-and-Merge Contour Models

We propose a new method for detection and tracking of moving objects from a moving camera image sequence using robust statistics and active contour models. First, we explain the outline of the proposed method in Sect. 2.1. We then describe the concrete algorithm for cancellation of apparent background motion based on LMedS estimation, and the method for tracking moving objects using the split-and-merge contour model with the energy function defined by the subtraction image in Sects. 2.2 and 2.3, respectively.

2.1 Outline of the Algorithm

In order to detect moving objects in a moving camera image sequence, we have to cancel the apparent motion of static background. We estimate the apparent background motion using LMedS (Least Median of Squares) method [10], [11]. After canceling the apparent background motion, we apply split-and-merge contour models [12], [13] for tracking multiple moving objects. Figure 1 shows an illustration of the proposed method. The method consists of the following four processes. Processes P1) to P3) are for the cancellation of apparent background motion to be explained in Sect. 2.2, and the process P4) is for tracking multiple moving objects to be described in Sect. 2.3.

P1) Background motion estimation using LMedS
Estimate the apparent background motion using LMedS method assuming that the static background is dominant in an image frame and that the apparent background motion between two consecutive image frames can be represented by affine transformation.

P2) Motion compensation
Cancel the apparent background motion by transforming the previous image using the affine parameters estimated in the process P1).

P3) Moving object detection
Detect moving objects by subtracting the affine transformed image in the process P2) from the current im-

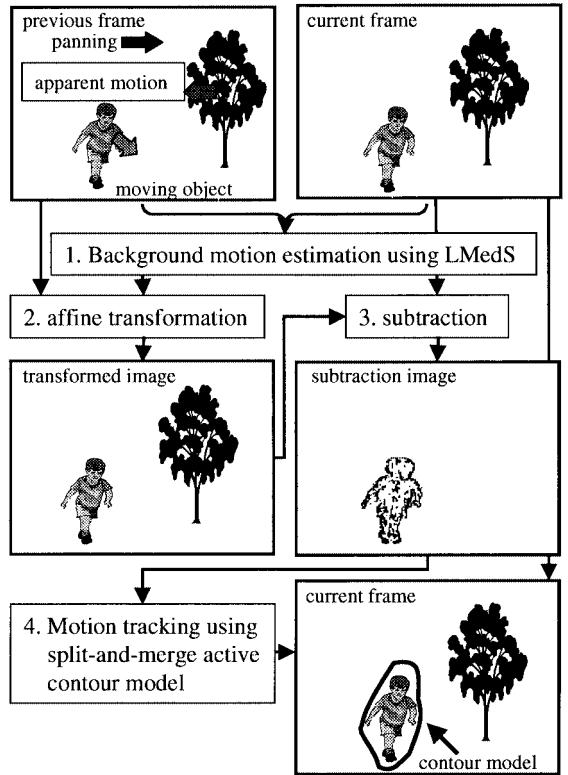


Fig. 1 Outline of the proposed method.

age.

P4) Contour tracking of moving objects

Track moving objects using the split-and-merge contour model whose energy function is defined based on the subtraction image in the process P3).

2.2 Cancellation of Apparent Background Motion Using LMedS Estimator

In this section, we explain a method of canceling the apparent background motion according to the processes P1) to P3) described in Sect. 2.1 in more detail.

First, we explain the process P1) for background motion estimation using LMedS. In this study, we assume that the static background is dominant in an image frame (i.e. the region of moving objects is less than 50% of the whole image) and that the background motion between two consecutive image frames can be represented by affine transformation. In general, when a camera moves for following moving objects, the change in the gaze direction between two consecutive frames is small. Under this constraint, when the change in the depth of the background is small or the distance between the eye point and the background is large enough compared to the distance of gaze points, possible transformation between two consecutive image frames is the combination of translation, rotation and scaling which

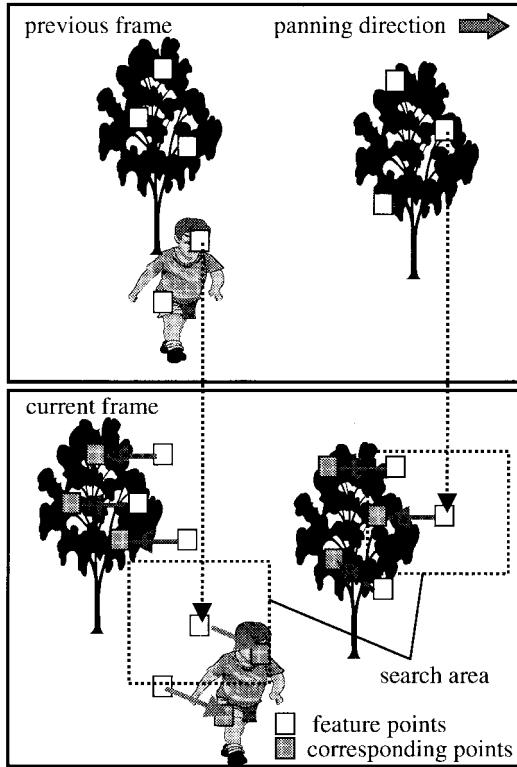


Fig. 2 Motion vector detection.

can be represented by affine transformation. Affine transformation parameters can be calculated using at least three feature points in the previous frame and their correspondences in the current frame. In order to estimate accurate background motion, we have to select feature points only from the background region and find their correct correspondences. We apply LMedS estimation to determine affine transformation parameters using correspondences for motion compensation of the background between two consecutive image frames according to the following procedure.

[STEP1]

Select N feature points in the previous frame at random where intensity variances of their small neighboring regions are large and search their corresponding points in the current frame using a standard correlation-based block matching technique. Figure 2 illustrates this step. White squares and meshed ones shown in Fig. 2 are feature points and their correspondences, respectively. A camera is assumed to be panning toward the right direction following a person walking to the right, so that apparent motion vectors of static background are toward the left. In this case, two feature points are selected in the moving object region and corresponding points are not always found correctly. Such points are outliers in the estimation of apparent background motion.

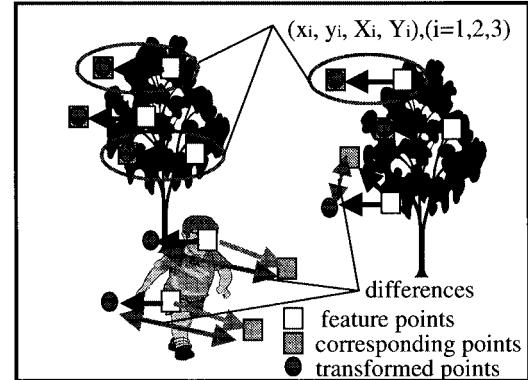


Fig. 3 Evaluation of affine parameters.

[STEP2]

Select a data set $(x_i, y_i, X_i, Y_i), (i = 1, 2, 3)$ at random (for example, it is indicated by elliptical marks shown in Fig. 3), which consists of three points $(x_i, y_i), (i = 1, 2, 3)$ from N feature points in the previous frame and their corresponding points $(X_i, Y_i), (i = 1, 2, 3)$ in the current frame. And then, calculate affine parameters from a_1 to a_6 in the following equation using the selected data set.

$$\begin{pmatrix} X_i \\ Y_i \end{pmatrix} = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} \begin{pmatrix} x_i \\ y_i \end{pmatrix} + \begin{pmatrix} a_5 \\ a_6 \end{pmatrix} \quad (1)$$

[STEP3]

Transform N feature points in the previous image using the parameters estimated in STEP2 and calculate a median of squared differences between the locations of transformed points and those of corresponding points searched in STEP1. For example, meshed circles shown in Fig. 3 denote the transformed points.

[STEP4]

Repeat STEP2 and STEP3 M times and select the estimated affine parameters for which the median of squared difference is the minimum.

According to the above procedure, the probability p that at least one data set of three points in the background and their correct corresponding points are obtained is derived from the following equation:

$$p(\epsilon, q, M) = 1 - \left(1 - ((1 - \epsilon)q)^3 \right)^M, \quad (2)$$

where $\epsilon (< 0.5)$ is the ratio of moving object regions to a whole image, q is the probability that corresponding points are correctly found in STEP1. For example, $p \approx 0.993$ when $\epsilon = 0.3$, $q = 0.7$ and $M = 40$.

Next, we explain the processes P2) for motion compensation and P3) for moving object detection described in Sect. 2.1 by using sample images. Moving object regions can be detected by subtracting the previous frame transformed with estimated affine param-

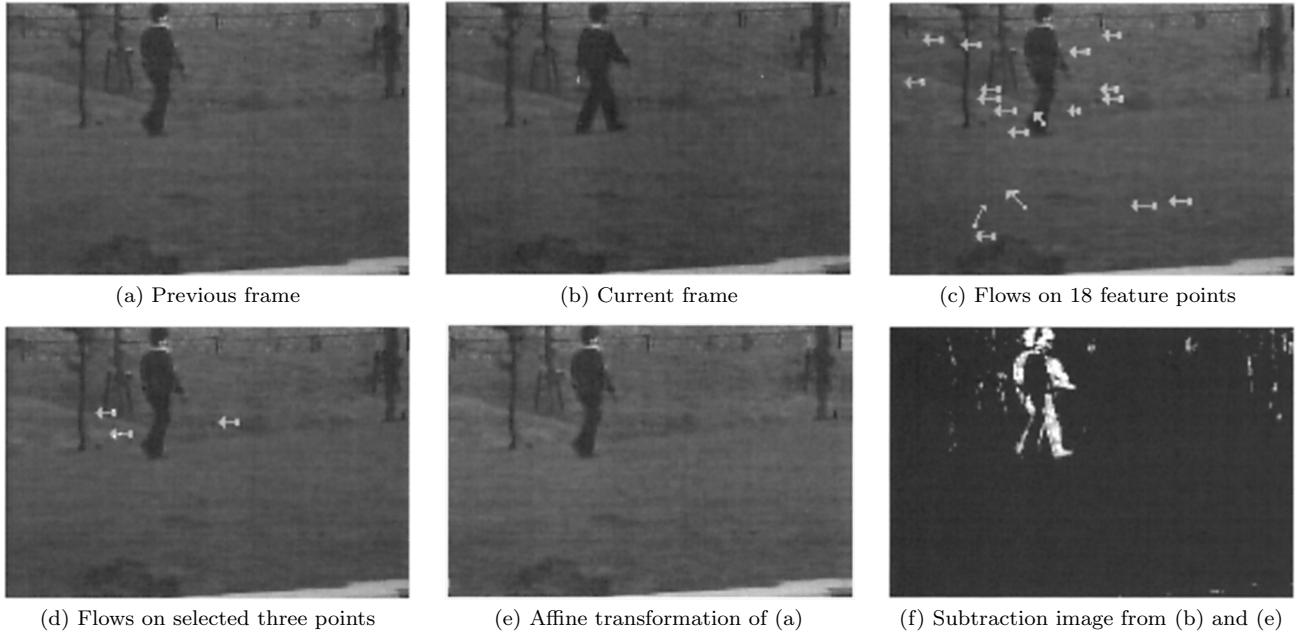


Fig. 4 Registration of a static background.

eters according to the process P1) from the current frame. Figures 4(a) and (b) are examples of two consecutive image frames. A camera is panning toward the right direction following the person walking to the right. As shown in Fig. 4(c), corresponding points are not always found correctly and one feature point is selected in the moving object region. Such points are outliers in the estimation of apparent background motion. Figure 4(d) shows the selected points and their flows which are used for calculating affine parameters. Figure 4(e) is the transformed image of the previous frame shown in Fig. 4(a) using the estimated affine parameters. Figure 4(f) shows the subtraction image whose pixel intensity is defined by the absolute value subtracted the transformed image shown in Fig. 4(e) from the current image shown in Fig. 4(b). A moving object region is almost correctly extracted, but there also exist extraction errors in the background which is mainly caused by approximation errors to a real scene using affine transformation and localization errors in corresponding point search in STEP1.

2.3 Tracking of Moving Objects by Split-and-Merge Contour Models

In this section, we explain a method of contour tracking of moving objects in process P4) described in Sect. 2.1 in more detail.

In the proposed method, we employ active contour models (snakes) [7] for tracking multiple moving objects. Subtraction images generated by the process described in Sect. 2.2 are utilized for defining an energy function of snakes. Snakes are parameterized curves $\mathbf{v}(s) = (x(s), y(s))$ ($0 \leq s \leq 1$) on an image plane

(x, y) , which are deformed to detect object boundaries by minimizing the following energy functional:

$$E_{\text{Snakes}} = E_{\text{int}}(\mathbf{v}) + E_{\text{image}}(\mathbf{v}) + E_{\text{ext}}(\mathbf{v}), \quad (3)$$

where E_{int} is an internal energy associated with splines, E_{image} is an image energy such as edge potential and E_{ext} is an external energy associated with external forces. In the present work, the image energy E_{image} is defined as

$$\begin{aligned} E_{\text{image}}(\mathbf{v}) &= E_{\text{edge}}(\mathbf{v}) = -\frac{1}{2}w_{\text{edge}}|\nabla I(\mathbf{v})|^2 \quad (4) \\ w_{\text{edge}} &= \begin{cases} \gamma; & I_s(\mathbf{v}) \geq T \\ 0; & \text{otherwise} \end{cases} \end{aligned}$$

where E_{edge} is the edge potential defined by a spatial gradient of intensity $I(\mathbf{v})$ of an original image, w_{edge} is a weight for E_{edge} , $I_s(\mathbf{v})$ is a pixel intensity of subtraction image after motion compensation by the method described in Sect. 2.2, γ is a predetermined weight and T is a predetermined threshold. From the definition above, the image energy is zero when a contour point is in the regions where intensity of the subtracted image is small so that contour models are attracted only to edges of moving objects.

Our goal is tracking of multiple moving objects in a moving camera image sequence, however, conventional snakes have the following two difficulties in their application to motion tracking: 1) automatic extraction of multiple objects in the first frame in a dynamic image, 2) tracking of apparently overlapping objects. The first problem primarily concerns initial setting of contour models. In most of the existing methods, initial contours for moving target objects have to be carefully set around the actual boundary of each object because

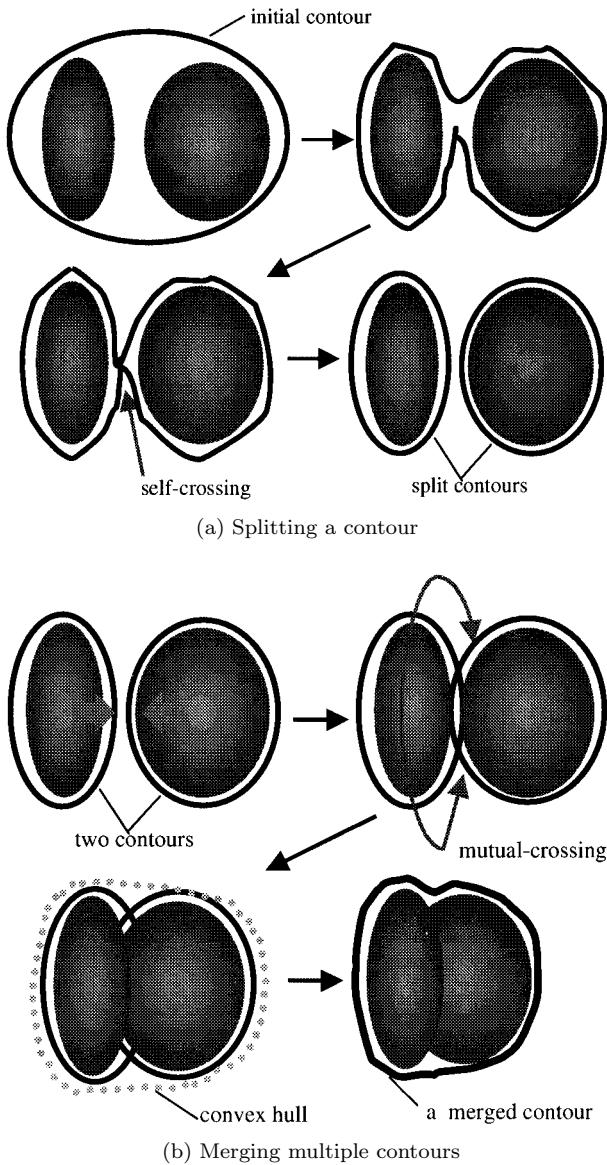


Fig. 5 An illustration of the split-and-merge contour model.

snakes including multiple objects can not extract each of them individually. The second problem in tracking multiple moving objects is that snakes can not track their boundaries correctly when they apparently overlap, since a contour model is attracted by edges which are extracted by another contour.

In order to overcome these difficulties, we have recently developed a split-and-merge contour model based on detecting self- and mutual-crossings of the contour and its feasibility have been proved in application to automatic tracking of multiple moving objects in a dynamic image from a static camera [12], [13]. We employed the split-and-merge contour model in the present work. In the method, an initial single contour is set on the first frame which is surrounding all moving objects (for example, it is simply selected as an image frame) and then iteratively splits into multiple

contours by cutting at self-crossing points for extracting multiple moving objects individually as illustrated in Fig. 5 (a). The method also detects mutual-crossing of different contour models which are extracting different objects for finding apparent overlapping of multiple moving objects. Multiple different contours which are mutually crossing are merged into a single contour as shown in Fig. 5 (b). A merged contour is defined as a convex hull of mutual-crossing contours. When overlapping objects leave away from each other, the merged contour is split into multiple contours again. In this case, it is necessary to correctly match objects before and after their overlapping. We employ a simple matching method based on computing a textural feature of an object region. In each frame, a cumulative intensity histogram is computed as the feature of an object.

By employing split-and-merge contour models with the image energy defined in Eq. (4), multiple moving objects can be extracted and tracked in a moving camera image sequence.

3. Real-Time Tracking of Moving Objects Using DSP Boards

In order to realize real-time tracking of multiple moving objects, we have implemented the method on an experimental image processing system mainly constructed of a video I/O board and eight DSP boards ISHTAR-II [14] containing TMS320C40.

3.1 System Configuration

The system configuration is illustrated in Fig. 6. The video board digitizes NTSC composite signals from a video camera into images which consist of 160×120 pixels. Each DSP board processes digitized images and transfers both the current result and the original image to the next board at every $1/7.5$ second according to a program which is down-loaded from a host computer (SUN SparcStation20) through the MPU board (S-BUS/VME converting interface). Installed programs in each DSP boards are as follows.

[DSP1-3]

Six feature points with large intensity variances of their small neighboring regions are selected in the previous frame at random in parallel in each DSP boards and their corresponding points are searched in the current frame. By using these three boards, the number of selected feature points in each frame is $N = 18$ in total. In the search process of corresponding points, block size and search area are 5×5 and 21×21 pixels, respectively. Each DSP board has two frame memories, the current frame from video board is stored in one memory and the previous frame is in the other. DSP1 to 3 send the locations of feature points and their correspondences to DSP4. DSP1 also sends the previous frame

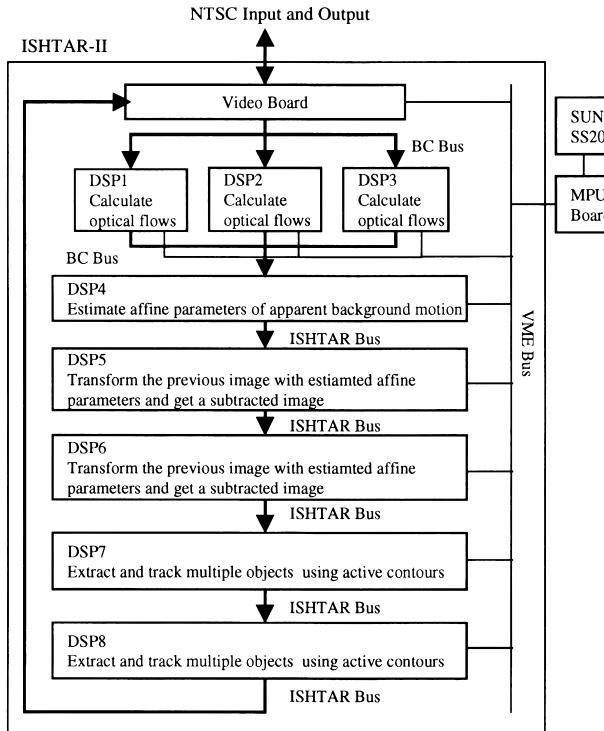


Fig. 6 System configuration of a DSP-based experimental platform.

to DSP4.

[DSP4]

DSP4 estimates affine parameters representing apparent background motion based on LMedS estimation described in Sect. 2.2. Provided that the ratio ϵ of moving object regions to a whole image is 0.3, the probability q that corresponding points are correctly searched is 0.7, and the number of iterations M in LMedS is set to be 40 times, the probability p that at least one data set of three feature points in the background and their correct correspondences are selected is 0.993. In this case, estimation error occurs only one frame per 20 seconds in average when DSP boards process 7.5 frames per second. DSP4 sends the estimated affine parameters and the image frame to DSP5.

[DSP5-6]

DSP5 and DSP6 subtract the previous frame transformed using estimated affine parameters from the current frame. In order to reduce the computational time, DSP5 and DSP6 process the upper half and lower half of the original image, respectively. DSP6 sends the subtracted image and the original current frame to DSP7.

[DSP7-8]

The split-and-merge contour model is implemented on DSP7 and DSP8 for tracking multiple moving objects. DSP7 receives the subtracted image and the original

image from DSP6. When the program loaded in DSP7 runs, an initial single contour is generated first on the first frame. The initial contour consists of 36 discrete snake points which are arranged in a rectangular shape along with the image frame. For every 1/7.5 second, DSP7 searches moving object boundaries five times in the current frame based on the pixel intensity of the subtracted image. Parameters γ and T of the image energy defined in Eq. (4) are set to 5.0 and 28 in this experiment, respectively. At any time when a new frame comes in, contour models are expanded for surrounding the moving objects in order to prevent that moving objects escape from contour models. After searching object boundaries five times, DSP7 sends the current number of contours, the current locations of discrete points of each contours, the subtracted image and the current image frame to DSP8. DSP8 also searches object boundaries four times for every 1/7.5 second by using a data set of contours and images transferred from DSP7. Since these two DSP boards are pipe-lined, object boundaries can be searched nine times in total in the same frame. DSP8 draws contour models on the current frame and sends the image to the video board for monitoring the tracking result.

3.2 Experimental Results

We show two experimental results of real-time tracking using the system described in Sect. 3.1.

Figure 7 shows an experimental result of real-time tracking of a walking person followed by a moving camera. Figure 7(a) shows three snapshots in a moving camera image sequence where a person is walking to the right direction in image frames. In each frame shown in Fig. 7(b), the region of walking person is extracted with some extraction errors in the background, which is mainly caused by approximation errors to a real scene using affine transformation and localization errors in searching corresponding points. Since extraction errors are small and discontinuous, the person is almost extracted and successfully tracked by the active contour model which smoothly connects pixels with motion cues as shown in Fig. 7(c). Tracking error in the small region near the persons feet is caused by his pivoting foot which does not move largely.

Figure 8 shows an experimental result of real-time tracking of two persons walking to the left direction in a moving camera image sequence. In earlier frames, the taller person is following the smaller person, and then, the taller person catches up and pass the smaller person in the succeeding frames. Figure 8(a) shows an initial contour which consists of 36 discrete points surrounding two persons. As shown in Figs. 8(b) and (c), the initial contour splits into two distinct contours which are detecting two persons, individually. In the current system, contour models can not converge to object boundaries at earlier image frames since search

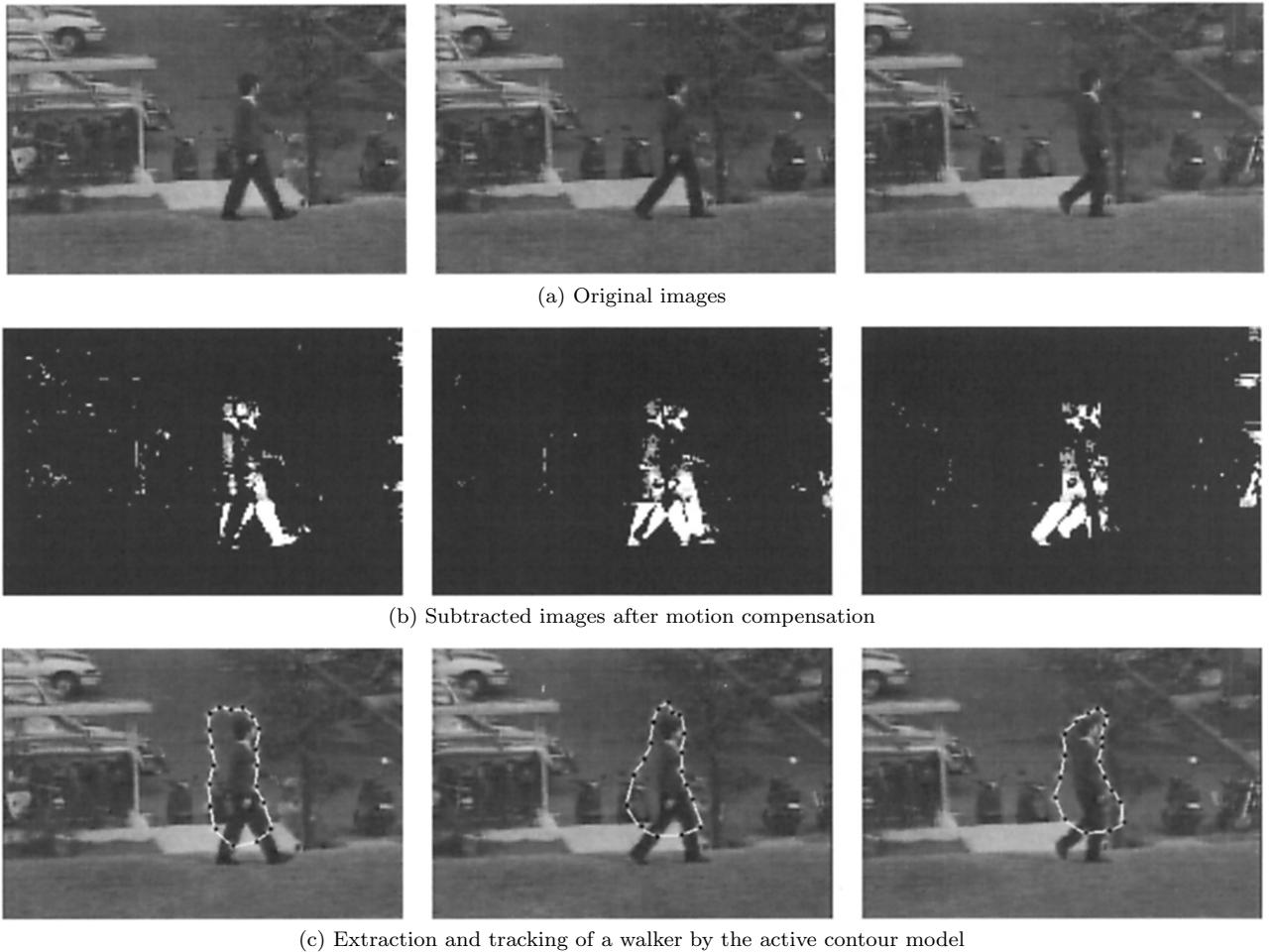


Fig. 7 Motion detection and tracking of a single moving object.

times of object boundaries can not be increased more than nine times under the limitation of processing speed of current DSP. Detected two persons are tracked in the succeeding frames shown in Figs. 8(d), (e) and (f). As shown in Fig. 8(g), two contours are merged into a single contour when the taller person catches up the smaller one. When the taller person is passing the smaller one, apparently overlapping two persons are tracked by the merged single contour shown in Figs. 8(h),(i) and (j). As two persons leave away from each other, the merged contour is deformed to occur the self-crossing (Figs. 8(k),(l),(m) and (n)), and finally it splits into two contours for tracking two persons again (Fig. 8(o)).

In these experiments, the time lag from input to output of an image frame is 0.8 second since six processes in which each image frame is processed in 1/7.5 second are pipe-lined.

4. Conclusions

In this paper, we have described a new method for real-time detection and tracking of multiple moving ob-

jects in a moving camera image sequence. The method is based on integrating robust estimation of apparent background motion and split-and-merge contour models. we have implemented the method on a DSP-based experimental image processing system and have shown the excellent performance for real-time tracking of multiple moving objects. Our future work includes the tracking of temporarily stopping objects in several frames since the present method can tracks only continuously moving edges.

Acknowledgement

This work was supported in part by Grant-In-Aid for Scientific Research under Grant Nos. 09221219 and 09555127 from the Ministry of Education, Science, Sports and Culture.

References

- [1] D. Murray and A. Basu, "Motion traking with an active camera," IEEE Trans. Pattern Anal. & Mach. Intell., vol.16, no.5, pp.449–459, May 1994.
- [2] D. Nair and J. Aggarwal, "Detecting unexpected moving

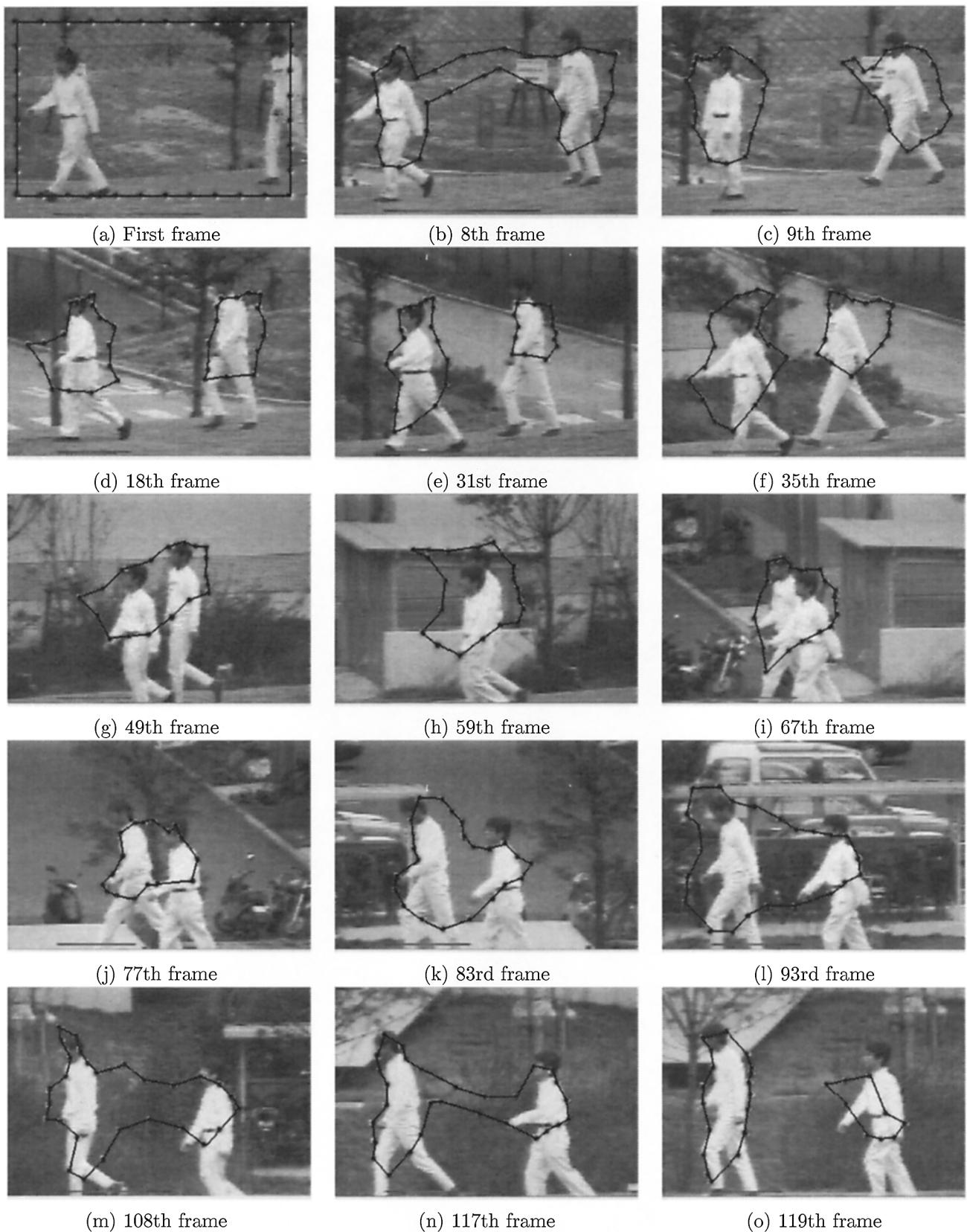


Fig. 8 Real-time tracking of multiple moving objects (two walkers) using DSP boards.

- obstacles that appear in the path of navigating robot," Proc. ICIP'94, vol.2, pp.311–315, 1994.
- [3] J. Odobezi and P. Bouthemy, "Detection of multiple moving objects using multiscale MRF with camera motion compensation," Proc. ICIP'94, vol.2, pp.257–261, 1994.
 - [4] M. Ben-Ezra, S. Peleg, and B. Rousso, "Motion segmentation using convergence properties," Proc. ARPA Image Understanding Workshop, pp.1233–1235, 1994.
 - [5] M. Iraniand, B. Rousso, and S. Peleg, "Detecting and tracking multiple moving objects using temporal integration," Proc. ECCV'92, pp.282–287, 1992.
 - [6] P. Rousseeuw and A. M. Leroy, Robust Regression & Outlier Detection, John Wiley & Sons, 1987.
 - [7] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active contour models," Proc. 1st ICCV, pp.259–268, 1987.
 - [8] H. Sawhney and S. Ayer, "Compact representations of videos through dominant and multiple motion estimation," IEEE Trans. Pattern Anal. & Mach. Intell., vol.18, no.8, pp.814–830, 1996.
 - [9] M. Etoh, "Promotion of block matching: Parametric representation for motion estimation," Proc. ICPR'98, vol.1, pp.282–285, 1998.
 - [10] T. Matsuoka, S. Araki, K. Yamazawa, H. Takemura, and N. Yokoya, "Tracking of moving objects in moving camera images and its real-time processing by DSP boards," IEICE Technical Report, PRMU97-235, 1998.
 - [11] S. Araki, T. Matsuoka, H. Takemura, and N. Yokoya, "Real-time tracking of multiple moving objects in moving camera image sequences using robust statistics," Proc. ICPR'98, vol.2, pp.1433–1435, 1998.
 - [12] S. Araki, N. Yokoya, and H. Takemura, "Tracking of multiple moving objects using split-and-merge contour models," Proc. Vision Interface '97, pp.65–72, 1997.
 - [13] S. Araki, N. Yokoya, and H. Takemura, "Real-time tracking of multiple moving objects using split-and-merge contour models based on crossing detection," IEICE Trans., vol.J80-D-II, no.11, pp.2940–2948, 1997.
 - [14] M. Shiohara, H. Egawa, S. Sasaki, M. Nagle, P. Sobey, and M. Srinivasan, "Real-time optical flow processor ISHTAR," Proc. 1st ACCV, pp.790–793, 1993.



Shoichi Araki received B.Eng. in Industrial Engineering from Osaka Prefectural University, Osaka, Japan, in 1989. He joined Matsushita Electric Industrial Co., LTD., in 1989, where he is currently a member of Central Research Laboratories. He was dispatched from the company to the Nara Institute of Science and Technology (NAIST) as a graduate student since 1993. He received M.Eng. and Ph.D. in Information Science from NAIST

in 1995 and 1998, respectively. Dr. Araki's research interests include computer vision and augmented reality.



Takashi Matsuoka received M.Eng. in Information Science from NAIST in 1998. He is currently a member of Mitsubishi Electric Corporation. His research interests include computer vision and robotics.



Naokazu Yokoya received the B.Eng., M.Eng and Ph.D. in information and computer sciences from Osaka University, Osaka, Japan, in 1974, 1976 and 1979, respectively. He worked at the Electrotechnical Laboratory of MITI from 1979 to 1993. Currently, he is a professor at the Nara Institute of Science and Technology, Ikoma, Nara, Japan. During the academic year 1986-1987, he was a visiting professor at the McGill Research Centre for Intelligent Machines, McGill University, Montreal, Quebec, Canada. Dr. Yokoya's research interests include image processing, computer vision and augmented reality. He is an associate editor of the journal of *Pattern Recognition*. He is a member of IEEE, IPSJ, VRSJ and JSAI, respectively.



Haruo Takemura received the B.Eng, M.Eng and Ph.D. in information and computer sciences from Osaka University, Osaka, Japan, in 1982, 1984 and 1987, respectively. He worked at the ATR Communication Systems Research Laboratories from 1987 to 1994. Currently, he is a associate professor at the Nara Institute of Science and Technology, Ikoma, Nara, Japan. During the academic year 1998, he was a visiting scholar at the Department of Mechanical and Industrial Engineering, University of Toronto, Toronto, Canada. Dr. Takemura's research interests include human-computer interaction, virtual reality and image processing. He is a member of IEEE, ACM, HFES, SPIE, IPSJ, HIS(J) and VRSJ, respectively.