

Simultaneous Tracking of Both Hands by Estimation of Erroneous Observations

James P. Mammen, Subhasis Chaudhuri and Tushar Agrawal
SPANN Lab, Department of Electrical Engg.,
Indian Institute of Technology, Bombay,
Powai, Mumbai-400076, INDIA.
sc@ee.iitb.ac.in

Abstract

The articulate motion of the hand makes it very difficult to track the hands while performing a gesture. Simultaneous tracking of both hands needs to deal with large interframe variations in shape, clutter and mutual occlusion. In this paper, we present a robust method for localizing the hand region by tracking even in the presence of severe occlusion. We develop a model for tracking two rectangular windows, each bounding one of the hands, using the condensation algorithm. We propose a new method for dealing with occlusions by estimating the occluded observations in terms of the non-occluded observations and their predicted values, yielding very reliable results.

1 Introduction

Hand tracking is an essential component of gesture recognition systems and haptic interfaces for virtual reality environments as mentioned in [1]. Sign language recognition and telerobotic applications also require a hand tracker. Multi-modal interfaces that provide natural human-computer interaction as well as a host of other applications such as robot programming by demonstration [2].

Various approaches to hand tracking have been used until now by researchers. Rehg and Kanade [3] describe a hand tracker using a 3D model of the hand with 27 degrees of freedom. The line features produced by fingers in grey-scale images without any background clutter are used as feature measurements. Goncalves *et al.* [4] use a 3D model of the arm with seven parameters. Using a single camera, they track the human arm against a dark uncluttered background with the help of a recursive estimator. Gavrila and Davis [5] use a 3D model with 22 degrees of freedom for the whole body with 4 degrees of freedom for each arm. They use multiple cameras to recover the 3D body pose at each time step. Obtaining the parameters of elaborate 3D models using computer vision techniques being quite complex, such systems have avoided issues of clutter in natural environments.

Deformable planar contours, known as snakes, coupled with the Kalman filtering framework can be used for tracking non-rigid objects like hands [6]. In [7], an improved version of the Kalman snake is used for tracking the hand which uses optical-flow to detect and reject image measurements corresponding to image clutter or other objects.

Isard and Blake [8] introduced a statistical factored sampling algorithm known as condensation for tracking in dense clutter. They apply it to track hand contours in a cluttered background. However, contour tracking techniques allow only a small subset of possible movements to maintain continuous deformation of contours. This limitation was overcome to some extent by Heap and Hogg [9], who describe an adaptation of the condensation algorithm for tracking across discontinuities in contour shapes. Possible shapes are represented as a union of a set of clusters in a high-dimensional shape space and discontinuous shape changes are described in terms of transitions between these clusters using a learned Markov model. However hand contours, being 2D projections of highly flexible human hands, assume too large a variation during natural motion of the hand to allow contour tracking.

Recently several approaches utilizing stereo ranging information have also been implemented. Lin [12] presents an algorithm for tracking the human arm using range image sequences. Jennings [13] describes a system for tracking the 3D position and orientation of a finger using several cameras. Tracking is based on combining multiple sources of information including stereo range images, color segmentation, shape information and various constraints.

Wren *et al.* [10] build a system which can track a single person by segmenting the image into blobs using color information and then using prior information about skin color and topology of a person's body to interpret the set of blobs as a figure. Bregler [11] describes a method for tracking human motion by grouping pixels with coherent motion, color and temporal support into blobs using an expectation-maximization (EM) algorithm and track each blob using a Kalman filter. These blob based methods are able to track even when there are greater variations from frame to frame. However, these methods do not deal with occlusions explicitly. MacCormik and Blake [14] track contours of two similar objects using an observation model which does not allow a single feature to correspond to both objects. In our work the observation model discourages the same feature from corresponding to both hands simultaneously in a probabilistic manner without disallowing it, since it is possible for both hands to be over each other temporarily. Utsumi and Ohya [15] track 3D positions and postures of both hands using multiple cameras. Each hand position is tracked with a Kalman filter and 3D hand postures are estimated using image features. This work deals with the mutual hand-to-hand occlusion inherent in tracking both hands, by selecting the camera images that do not have such occlusions.

The proposed work is a part of an ongoing project for gesture recognition to be used for telerobotic applications. Natural hand gestures display large interframe variations in position, posture and contour and hence we do not track the 3D pose or the contour of the hand. We track a rectangular window of variable size surrounding the hand region using the skin color as our observations. Thus, in effect, we track skin-colored blobs. The evolution of the size and the trajectory defined by the window is subsequently used in recognizing the gesture performed. This approach, as opposed to other methods, allows tracking of relatively faster motion with almost arbitrary variation in hand shape from frame to frame. The information obtained this way is satisfactory for most applications like gesture recognition. Moreover, it is possible to derive more exact information for applications which require so, by processing this localized region. Since we use skin colored blobs we need to deal with clutter due to possible presence of skin colored objects in the background and the face region. Hand-hand occlusion is a major problem in simultaneous tracking of both hands. We propose a robust method of tracking in the presence

of occlusion also by estimating the occluded observations. We use the condensation algorithm for tracking since this allows us to maintain multiple hypotheses in the observation model and thus reduces the possibility of loss of tracking.

The organization of the paper is as follows. The next section gives an overview of the condensation algorithm that we have used. This is followed by the stochastic dynamical models for the states and the observations for hand tracking. Thereafter, we present a robust method for dealing with occlusions and clutter. Finally experimental results of applying this method are presented.

2 Condensation Algorithm

Blake and Isard [8] proposed a stochastic algorithm for tracking curves in clutter using conditional density propagation over time, leading to the name condensation. It is a factored sampling approach to propagate the entire probability distribution of the parameters to be tracked over time. A summary of the algorithm and its framework follows.

The state of the blob-window to be tracked is denoted by X_t and its history by $\mathcal{X}_t = \{X_1, \dots, X_t\}$. The observation at time t is denoted by Z_t and its history by $\mathcal{Z}_t = \{Z_1, \dots, Z_t\}$. It is assumed that the state dynamics form a temporal Markov chain so that $p(X_t|\mathcal{X}_{t-1}) = p(X_t|X_{t-1})$. Moreover, the measurements are assumed to be independent, both mutually and with respect to the dynamical process. Under these assumptions, it is shown that the conditional state-density at time t is given by

$$p(X_t|\mathcal{Z}_t) = k_t p(Z_t|X_t) p(X_t|Z_{t-1}) \quad (1)$$

where k_t is a normalization constant that does not depend on X_t . In the condensation algorithm, the conditional state density is approximated by a sample-set of size N , $S_t = \{s_t^{(1)}, \dots, s_t^{(N)}\}$ and a corresponding set of weights $\Pi_t = \{\pi_t^{(1)}, \dots, \pi_t^{(N)}\}$, each given by

$$\pi_t^{(i)} = \frac{p(Z_t|s_t^{(i)})}{\sum_{j=1}^N p(Z_t|s_t^{(j)})} \quad (2)$$

The condensation algorithm constructs from the sample-set at time $t-1$, i.e. $\{S_{t-1}, \Pi_{t-1}\}$, a new sample-set for time t , i.e. $\{S_t, \Pi_t\}$. The n^{th} of the N samples is constructed as follows:

1. Select a sample $s'_t = s_{t-1}^{(j)}$ with probability $\pi_{t-1}^{(j)}$.
2. Predict by sampling from $p(X_t|X_{t-1} = s'_t)$ to choose each $s_t^{(n)}$.
3. Obtain the new weight given the observation Z_t using $\pi_t^{(n)} = p(Z_t|X_t = s_t^{(n)})$.

After all the N weights are calculated they are normalized to obtain the new set of weights.

An estimate of any moment of the tracked position at time t is given by,

$$E[q(X_t)] = \sum_{n=1}^N \pi_t^{(n)} q(s_t^{(n)}) \quad (3)$$

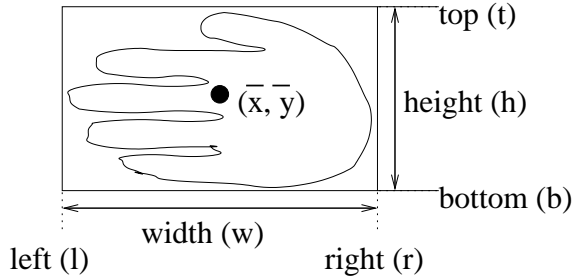


Figure 1: Illustration of the state defined by the window bounding the hand (palm).

3 Tracking a Single Hand

To develop a hand tracker using the condensation algorithm we develop a model for the state dynamics and a measurement model.

3.1 Model for State Dynamics

Since our purpose is to track rectangular windows bounding the hands, we select the coordinates of the center of each rectangular window and its height and width, as shown in Fig. 1, as elements of the 4-dimensional state vector $x_t = [\bar{x} \ \bar{y} \ h \ w]^T$. We model the state dynamics as a second-order AR process.

$$x_t = A_2 x_{t-2} + A_1 x_{t-1} + w_t \quad (4)$$

where w_t is a zero-mean, white Gaussian random vector. This is intuitively satisfying, since the state dynamics may be thought of as two dimensional translation and change in size of a rectangular window surrounding the hand region. We form an augmented state vector X_t for each window as follows.

$$X_t = \begin{pmatrix} x_{t-1} \\ x_t \end{pmatrix} \quad (5)$$

Thus we may rewrite eqn. (4) as $X_t = AX_{t-1} + W_t$. The Markovian nature of this model is evident in above equation.

3.2 Observation Model

In order to differentiate the hand region from the rest of the image we need a strong feature which is specific to the hand. It has been found that irrespective of race, skin color occupies a small portion of the color space [16]. As a result, skin color is a powerful cue in locating the unadorned hand.

In order to detect skin colored regions we assign to each pixel at location Y , a number $n(Y)$ which is proportional to the ratio of probabilities of its being skin to its not being skin.

$$n(Y) \propto \frac{P(\text{skin}|C_Y)}{P(\text{notskin}|C_Y)} \quad (6)$$

where $C_Y = [Cb_Y, Cr_Y]^T$ is the color of the Y^{th} pixel in the $CbCr$ space. We use the $YCbCr$ color space [17] representation and use only the Cb and Cr values to eliminate the effect of intensity variation. Using Bayes theorem, the above equation can be rewritten as

$$n(Y) \propto \frac{P(C_Y|skin)P(skin)}{P(C_Y|notskin)P(notskin)} \quad (7)$$

Since $P(skin)$ and $P(notskin)$ are constants irrespective of the color of the pixel, they can be eliminated from above equation.

The likelihood functions $P(C_Y|skin)$ and $P(C_Y|notskin)$ are obtained by learning from a large number of images. Based on a histogram of $n(X)$, we select pixels having very high values of $n(X)$ as seeds and starting from them, we form skin colored blobs by including connected pixels having the likelihood ratio above an appropriate lower threshold which is also based on the histogram. Thus we avoid a simple thresholding scheme by using the histogram of $n(X)$ and also impose a connectivity constraint resulting in skin colored blobs. It should be noted that using skin color detection will yield regions not only of the hands but also of the face and the neck. Apart from this, even other objects like wood are likely to be classified as skin.

The observations are the top, left, right and bottom edges delimiting the blobs as shown in the Fig. 1. Thus the observation vector at time t is given by $Z_t = [t \ l \ b \ r]^T$. There may be one or more possible measurements Z_t^1, \dots, Z_t^m corresponding to m skin-colored blobs, with more than one blob being detected due to clutter. Having selected the state as explained in the previous subsection, we model the true observation to arise from the state as

$$Z_t = BX_t + F_t \quad (8)$$

where F_t is assumed to be a zero mean, white Gaussian random noise. Thus in the absence of clutter we could write $P(Z_t|X_t)$ in (1) as $p(Z_t|X_t) \propto \exp -(Z_t - BX_t)^T R_F^{-1} (Z_t - BX_t)$ where R_F is the covariance matrix of the random vector F_t .

However, in order to deal with multiple observations, arising due to more than one blob being detected owing to clutter, we use the following observation model.

$$p(Z_t|X_t) \propto \exp\{-(Z_t^j - BX_t)^T R_F^{-1} (Z_t^j - BX_t)\} \quad (9)$$

where $j = \arg \min_i (Z_t^i - BX_t)^T R_F^{-1} (Z_t^i - BX_t)$.

In effect, for each given state sample, we select the measurement corresponding to the blob which has the maximum probability of being the measurement. Most of the times only one skin colored blob will be detected. Only in the case of clutter does more than one blob arise. Hence, for a practical implementation of the hand tracker we select the blob which fits the predicted value $BX_{t|t-1}$ best as the measurement, and use the same measurement for calculating weights of all samples as given by (2).

4 Models for Tracking Both Hands Simultaneously

For joint tracking of both hands we proceed as follows. For each hand a search window is formed around the predicted values and if the search windows have a significant overlap, the tracking is done jointly or else each hand is tracked independently as mentioned in the

previous section. After detecting skin-colored blobs, the occluded components, if any, of the observed blob are estimated, as explained in the next section. When both the search windows overlap, there is a possibility that both the tracked blobs may latch on to the same blob. To avoid this, for joint tracking, we concatenate the state and observation vectors of the left and the right hand (indicated by the superscripts L and R), into X'_t and Z'_t respectively as follows

$$X'_t = \begin{pmatrix} X_t^L \\ X_t^R \end{pmatrix} \quad Z'_t = \begin{pmatrix} Z_t^L \\ Z_t^R \end{pmatrix}$$

and use $p(Z'_t|X'_t) \propto \exp\{-(Z_t^k - BX_t^L)^T R_F^{-1} (Z_t^k - BX_t^L) - (Z_t^l - BX_t^R)^T R_F^{-1} (Z_t^l - BX_t^R)\}$ as the observation model, where, $(k, l) = \arg \min_{i,j} (Z_t^i - BX_t^L)^T R_F^{-1} (Z_t^i - BX_t^L) + (Z_t^j - BX_t^R)^T R_F^{-1} (Z_t^j - BX_t^R) + K\delta(i, j)$.

This is similar to (9) except for the term $K\delta(i, j)$, which prevents both the windows from tracking the same blob by serving as a penalty constant. The effect of this term is similar to that of the method used for exclusion in [14] by not allowing the same feature to correspond to two different objects simultaneously. However, the term $K\delta(i, j)$ discourages a blob from being associated with both hands in a probabilistic manner without disallowing it. This is necessary since occasionally both the hands may overlap to form a single blob. The value of K should be chosen appropriately based on experimentation. Again, for a practical implementation we select the pair of blobs which fits the predicted value best to be the observation and allow multiple observations only when more than one pair of blobs closely fits the predicted value.

5 Providing Robustness to the Tracker

As mentioned in the previous section, we form a search window based on the predicted values and detect skin-colored blobs. However the observations may be cluttered due to the face, objects having skin-like color or false skin color detection. While tracking both hands, the observations are in error when the hands obstruct each other. However, usually all the elements of the observation vector are not simultaneously in error. Hence we propose a scheme for estimating the erroneous observations based on the ones that are not erroneous and the predicted values of the states. We need a method for determining whether an observation is erroneous or not. If any observed value does not match the predicted value within a threshold, we assume it to be due to occlusion or clutter. Thus the observations could be erroneous due to the mutual obstruction of hands, background objects similar to skin, body parts like face and throat or false skin color detection or even error in prediction.

Taking this into consideration, we may rewrite the observation model in (8) as

$$Z_t = BX_t + D_t F_t + D_t^c Q_t \quad (10)$$

where, F_t is the observation noise as in (8) and Q_t corresponds to the large perturbations in the measurement due to occlusion or clutter. Let the observation vector Z_t be of dimension $M \times 1$ and the state X_t be of dimension $K \times 1$ resulting in the dimension of matrix B being $M \times K$ and that of D_t and D_t^c being $M \times M$. All the non-diagonal entries of the matrix D_t are zero and the diagonal entries are either 0 or 1 depending on whether

the observation is erroneous or not. D_t^c also has all non-diagonal entries to be zero and the diagonal elements are the 1's complement of the corresponding elements in D_t . Consider the case when only the i^{th} observation is in error due to clutter or occlusion and the remaining measurements are not affected. Thus, given that the i^{th} element z_i in the observation vector is erroneous, we need to estimate z_i in terms of the other measurements and the predicted states. Substituting the measurements in (10) and taking expectation of both sides assuming that Q_t has zero mean, we get, $[z_1 \dots E(z_i) \dots z_M]^T = BE(X_t)$. Multiplying both sides with a row vector G , with elements $[g_1 \dots g_M]^T$, and rearranging we get, $g_i E(z_i) = \sum_{j \neq i} g_j z_j + GBE(X_t)$

Now we substitute $E(X_t)$ by the predicted values $\hat{X}_{t|t-1}$. Also letting $GB = H = [h_1 \dots h_K]$ and $g_i = 1$, we obtain,

$$E(z_i) = \sum_{j \neq i} g_j z_j + [h_1 \dots h_K] \hat{X}_{t|t-1} \quad (11)$$

Since we wish to estimate the erroneous observation in terms of the remaining measurements rather than depend on the predictions, the vector G should be selected such that the maximum number of elements in H are zero.

However, as mentioned earlier we consider an observation to be erroneous when it differs from the prediction by a value greater than some threshold. This conclusion could also be due to an error in the prediction and hence we take a linear combination of the estimate $E(z_i)$ obtained in (11) and the measured value z_i .

$$\hat{z}_i = pE(z_i) + (1 - p)z_i \quad 0 < p < 1 \quad (12)$$

The value of p should be selected depending on the confidence that the measurement is really erroneous and that prediction error is not the cause. This value can be decided based on the number of elements that are erroneous. If all the inter-related observables are erroneous, the estimate is based only on the predicted values. A lower value of p would be suitable in this case, since simultaneous high error in all the inter-related observations could be due to error in prediction itself. In the case of two hand tracking, when occlusion due to the other hand can be predicted, if occlusion occurs as predicted then the p can be made very close to 1. Our choice of the state and observation vectors makes t and b a pair of inter-related observations through the predicted states. Similarly l and r form another pair. Hence when only any one of either t or b is occluded we take $p = 0.8$ and similarly for the pair l and r . When both t and b or l and r are occluded we take $p = 0.6$. In the case where the occlusion occurs as predicted we take $p = 0.95$.

The estimated value obtained thus is used instead of the erroneous value and the tracking proceeds as before. In case, more than one observable is erroneous, the same method can be extended. If z_i and z_j are erroneous then the estimate of z_i should not include z_j and hence we need to make the term $h_j = 0$.

6 Experimental Results

The results of hand tracking for some of the frames in a video sequence captured at 12 frames per second using the model in section 3 in conjunction with the robustness technique given in section 5 are shown in Fig. 2. As can be seen, the hand is tracked in a

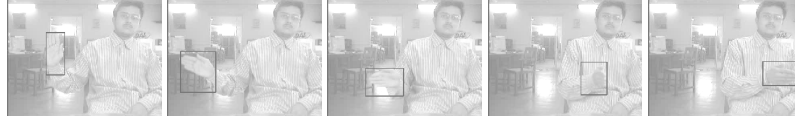


Figure 2: Results of tracking a single hand for the “Move Left” gesture.

usual office environment with background clutter. We run the condensation algorithm using just 100 sample points for the distribution allowing a fast implementation. The initial position of the hand is specified using a palm detection algorithm described in [18], [19]. The results show that the tracked window size increases beyond the hand due to false skin color detection but the method recovers later.

We show the results of tracking natural movements of both hands simultaneously in video sequences captured at 12 frames per second using the model in section 4 along with the robustness technique of section 5. The topmost rows of Figs. 3 and 4 show intermediate frames of two of the test sequences (complete video for these results are available on the CDROM in MPEG format). The two tracked windows are shown by rectangular boxes (one continuous and the other dotted) surrounding the hand. A region containing occlusion is shown frame by frame in the second row of each figure for greater details. The third row shows the observations obtained by detecting skin-colored blobs within the search windows indicated by the dotted rectangles. The presence of two search windows indicates that both hands are being tracked independently, whereas, a single search window indicates that joint tracking of both hands is being done due to significant overlap in the search windows as mentioned in section 4. As seen in Fig. 3, both the blobs merge for 2 frames and for the sequence shown in Fig. 4, the blobs remain merged for 5 frames. In spite of this, the tracker successfully follows each hand as shown by the two rectangular windows due to the estimation of occluded observations as in section 5. As a result of using the proposed observation model, the tracked windows separate even after the blobs remaining merged for 5 frames, as shown in Fig. 4.

7 Conclusion

We present a robust hand tracker which simultaneously tracks rectangular windows surrounding both the hands using the condensation algorithm. The tracking algorithm is able to deal with error in observations due to clutter and occlusions by estimating the occluded observations in terms of the unoccluded measurements and the predicted state values by exploiting the fact that usually all measurements are not occluded simultaneously. Results of applying this method to track both hands simultaneously show that this method is capable of dealing successfully with various kinds of clutter and hand-to-hand occlusion. The results of the tracker will, in future, be used in recognizing the performed gesture.

References

- [1] V. Pavlovic, R. Sharma and T. S. Huang, “Visual Interpretation of Hand Gestures for Human-Computer Interaction: A Review”, *IEEE trans. on PAMI*, 19(7), pp. 677-695, 1997.

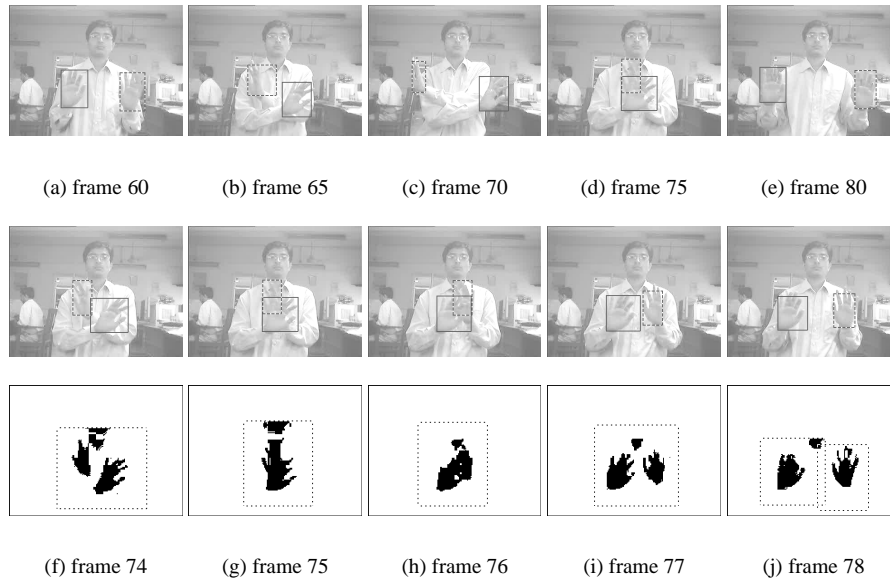


Figure 3: 1st row: Intermediate frames from a video sequence in which a two-handed ‘all clear’ gesture is performed. 2nd row: Frames during an occlusion phase. 3rd row: Corresponding skin-colored blobs. The dotted lines indicate the search windows. Each frame is of size 384×288 .

- [2] M. Yeasin and S. Chaudhuri, “Automatic Generation of Robot Program Code : Learning from Perceptual Data”, *In Proc. of Int. Conf. Computer Vision*, pp. 889-894, Mumbai, India, 1998.
- [3] J. Reh and T. Kanade, “Visual Tracking of High Dof Articulated 0: An Application to Human Hand Tracking”, *In Proc. of 3rd European Conf. on Computer Vision*, 2, pp. 35-46, Sweden, 1994.
- [4] L. Goncalves, E. Di Bernardo, E. Ursella and P. Perona, “Monocular Tracking of the Human Arm in 3D”, *In Proc. of IEEE Int. Conf. Computer Vision*, Cambridge, 1995.
- [5] D.M. Gavrila and L.S. Davis, “3-D model-based tracking of humans in action: a multi-view approach”, *In Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, San Francisco, 1996.
- [6] D. Terzopoulos and R. Szeliski, “Tracking with Kalman Snakes”, *Active Vision* ed. A. Blake and Yuille, 3-20, MIT, 1992.
- [7] Natan Peterfreund, “Robust Tracking of Position and Velocity With Kalman Snakes”, *IEEE trans. on PAMI*, 10(6), pp. 564-569, June 1999.
- [8] M. Isard and A. Blake “Condensation - Conditional Density Propagation For Visual Tracking,” *Int. Journal. of Computer Vision* , Vol. 28,1, pp. 5-28, 1998.
- [9] Tony Heap and David Hogg, “Wormholes in Shape Space: Tracking through Discontinuous Changes in Shape”, *In Proc. of ICCV*, pp. 344-349, Mumbai, 1998.
- [10] C. Wren, A. Azarbayejani, T. Darell, A. Pentland, “Pfinder: Real-Time Tracking of the Human Body”, *IEEE trans. on PAMI*, 19(7), 780-785, 1997.

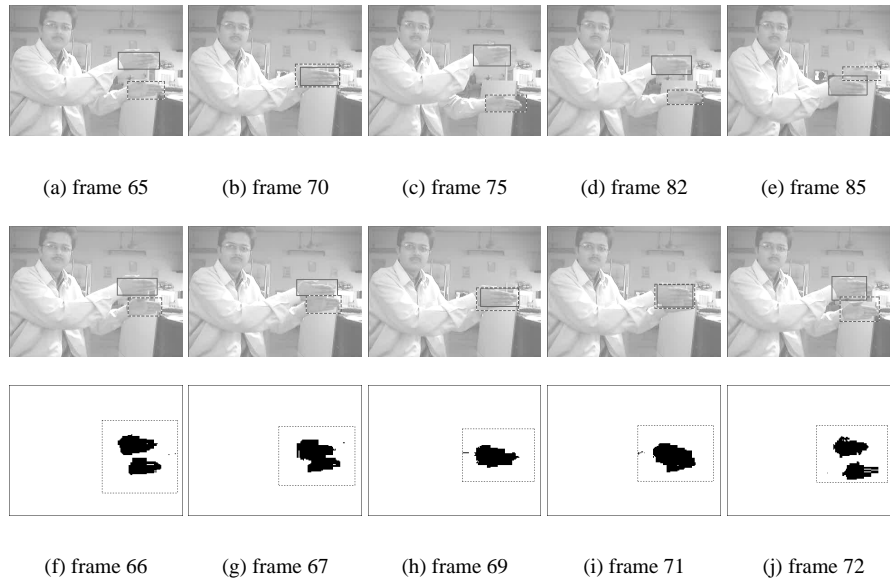


Figure 4: 1st row: Intermediate frames from a video sequence where a ‘chop it’ gesture is performed. The hands come closer, overlap and return back and thereafter cross each other. 2nd row: Frames during an occlusion phase. 3rd row: Corresponding skin-colored blobs.

[11] Christoph Bregler, “Learning and Recognizing Human Dynamics in Video Sequences”, *IEEE Conf. on Computer Vision and Pattern Recognition*, Puerto Rico, 1997.

[12] Michael Lin, “Tracking Articulated Objects In Real-Time Range Image Sequences”, *In Proc. of ICCV*, pp. 648-653, Corfu, Sep 1999.

[13] C. Jennings, “Robust finger tracking with multiple cameras”, *IEEE workshop on Recognition, Analysis and Tracking of Faces and Gestures in Real-Time Systems*, Corfu, 1999.

[14] J. MacCormick and A. Blake, “A Probabilistic Exclusion Principle for Tracking Multiple Objects”, *In Proc. of Int. Conf. on Computer Vision*, Corfu, Sep 1999.

[15] Akira Utsumi and Jun Ohya, “Multiple-Hand-Gesture Tracking using Multiple Cameras”, *In Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, 1, pp. 473-478, Colorado, 1999.

[16] R. Kjeldsen and J. Kender, “Finding skin in color images,” *In Proc. of the Second Int. Conf. on Automatic Face and Gesture Recognition*, pp. 312-317, 1996.

[17] Keith Jack, “Video Demystified: A Handbook for the Digital Engineer”, *HighText publications*, California, 1996.

[18] James Mammen and S. Chaudhuri, “A Model Based Technique for Palm Detection”, *In Proc. of the National Conf. on Communication*, pp. 315-319, Delhi, Jan 2001.

[19] James P. Mammen, “Hand Tracking and Gesture Recognition”, M. Tech. Thesis, Indian Inst. of Tech., Bombay, Jan 2001.