

Flux Tensor Constrained Geodesic Active Contours with Sensor Fusion for Persistent Object Tracking

Filiz Bunyak, Kannappan Palaniappan, Sumit Kumar Nath

Department of Computer Science, University of Missouri-Columbia, MO 65211-2060, USA

Email: {bunyak,palaniappan}@missouri.edu, naths@ecse.rpi.edu

Gunasekaran Seetharaman

Dept of Electrical and Computer Engineering, Air Force Institute of Technology, OH 45433-7765 USA

Email:guna.seetharaman@afit.edu

Abstract—This paper makes new contributions in motion detection, object segmentation and trajectory estimation to create a successful object tracking system. A new efficient motion detection algorithm referred to as the flux tensor is used to detect moving objects in infrared video without requiring background modeling or contour extraction. The flux tensor-based motion detector when applied to infrared video is more accurate than thresholding "hot-spots", and is insensitive to shadows as well as illumination changes in the visible channel. In real world monitoring tasks fusing scene information from multiple sensors and sources is a useful core mechanism to deal with complex scenes, lighting conditions and environmental variables. The object segmentation algorithm uses level set-based geodesic active contour evolution that incorporates the fusion of visible color and infrared edge informations in a novel manner. Touching or overlapping objects are further refined during the segmentation process using an appropriate shape-based model. Multiple object tracking using correspondence graphs is extended to handle groups of objects and occlusion events by Kalman filter-based cluster trajectory analysis and watershed segmentation. The proposed object tracking algorithm was successfully tested on several difficult outdoor multispectral videos from stationary sensors and is not confounded by shadows or illumination variations.

Index Terms—Flux tensor, sensor fusion, object tracking, active contours, level set, infrared images.

I. INTRODUCTION

In real world monitoring tasks, persistent moving object detection and tracking remains a challenging problem due to complex scenes, lighting conditions, environmental variables (particularly in outdoor settings with weather), clutter, noise, and occlusions. Fusing scene information from multiple sensors and sources is a useful core mechanism to deal with these problems. Recent developments in micro-optics and micro-electromechanical systems (MEMS), VCSELS, tunable RCLEDs (resonant cavity LEDs), and tunable micro-bolometers indicate that hyperspectral imaging will rapidly become as ubiquitous as visible and thermal videos are today [1], [2].

This paper is based on "Geodesic Active Contour Based Fusion of Visible and Infrared Video for Persistent Object Tracking," by F. Bunyak, K. Palaniappan, S. Nath and G. Seetharaman, which appeared in the Proceedings of the 8th IEEE Workshop on Applications of Computer Vision (WACV 2007), Texas, USA, February 2007. © 2007 IEEE.

On board lidars and radars have been used successfully in unmanned autonomous vehicles, extending their versatility well beyond what was demonstrated in the 1990's based on dynamic scene analysis of visible video only. Most of the autonomous vehicles competing in the recent DARPA Grand Challenge events used one or more lidar sensors to augment the video imagery, demonstrating intelligent navigation using fusion of multiple information sources [3]. Autonomous navigation in city traffic with weather, signals, vehicles, pedestrians, and construction will be more challenging.

Effective performance in persistent tracking of people and objects for navigation, surveillance, or forensic behavior analysis applications require robust capabilities that are scalable to changing environmental conditions and external constraints (ie visibility, camouflage, contraband, security, etc.) [4]. For example, monitoring the barrier around sensitive facilities such as chemical or nuclear plants will require using multiple sensors in addition to a network of (visible) video cameras. Both infrared cameras and laser-scanner based lidar have been used to successfully enhance the overall effectiveness of such systems. In crowds or busy traffic areas even though it may be impractical to monitor and track each person individually, information fusion that characterizes objects of interest can significantly improve throughput. Airport surveillance systems using high resolution infrared/thermal video of people can extract invisible biometric signatures to characterize individuals or tight groups, and use these *short-term* multispectral blob signatures to resolve cluttered regions in difficult video segmentation tasks.

This paper presents a new moving object detection and tracking system for surveillance applications using fusion of visible and infrared information. A preliminary version of the paper has been published in [5]. Infrared imagery is less sensitive to illumination related problems such as uneven lighting, moving cast shadows or sudden illumination changes (i.e. cloud movements) that cause false detections, missed objects, shape deformations, false merges etc. in visible imagery. But use of infrared imagery alone often results in poor performance since generally these sensors produce imagery with low signal-to-

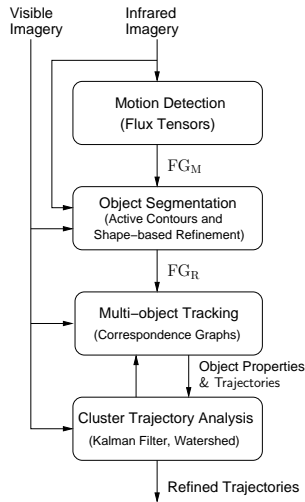


Figure 1: Multi-spectral data fusion system for persistent object tracking.

noise ratio, uncalibrated white-black polarity changes, and "halo effect" around hot or cold objects [6]. "Hot spot" techniques that detect moving objects by identifying bright regions in infrared imagery are inadequate in the general case, because the assumption that the objects of interest, people and moving cars are much hotter than the surrounding is not always true. The proposed system illustrated in Figure 1 is summarized below. A new efficient *motion detection* algorithm referred to as the flux tensor is used to detect moving objects in infrared video without requiring background modeling or contour extraction. The *object segmentation* algorithm uses level set-based geodesic active contour evolution that incorporates the fusion of visible color and infrared edge information. Touching or overlapping objects are further refined during the segmentation process using an appropriate shape-based model. *Multiple object tracking* module resolves frame-to-frame object correspondences, *cluster trajectory analysis* extension handles groups of objects and occlusion events.

This paper is organized as follows. In Section II motion detection using a flux tensor framework is explored. In Section III motion constrained object segmentation using geodesic active contours is discussed. In Section IV edge based video sensor fusion within a level set based active contour framework is explored. In Section V shape-based refinement of the object masks is presented. In Section VI multiple object tracking using correspondence graphs and its extension to handle groups of objects and occlusion events are described. Section VII presents the results, Section VIII offers the concluding remarks.

II. MOTION DETECTION USING A FLUX TENSOR FRAMEWORK

Motion blob detection is performed using our novel *flux tensor* method which is an extension to 3D grayscale structure tensor. Both the grayscale structure tensor and the proposed flux tensor use spatio-temporal consistency

more efficiently, thus produce less noisy and more spatially coherent motion segmentation results compared to classical optical flow methods [7]. The flux tensor is more efficient in comparison to the 3D grayscale structure tensor since motion information is more directly incorporated in the flux calculation which is less expensive than computing eigenvalue decompositions as with the 3D grayscale structure tensor.

A. 3D Structure Tensors

Structure tensors are a matrix representation of partial derivative information. As they allow both orientation estimation and image structure analysis they have many applications in image processing and computer vision. 2D structure tensors have been widely used in edge/corner detection and texture analysis, 3D structure tensors have been used in low-level motion estimation and segmentation [7], [8].

Under the constant illumination model, the optic-flow (OF) equation of a spatiotemporal image volume $\mathbf{I}(\mathbf{x})$ centered at location $\mathbf{x} = [x, y, t]$ is given by Eq. 1 [9] where, $\mathbf{v}(\mathbf{x}) = [v_x, v_y, v_t]$ is the optic-flow vector at \mathbf{x} ,

$$\begin{aligned} \frac{d\mathbf{I}(\mathbf{x})}{dt} &= \frac{\partial\mathbf{I}(\mathbf{x})}{\partial x} v_x + \frac{\partial\mathbf{I}(\mathbf{x})}{\partial y} v_y + \frac{\partial\mathbf{I}(\mathbf{x})}{\partial t} v_t \\ &= \nabla\mathbf{I}^T(\mathbf{x}) \mathbf{v}(\mathbf{x}) = 0 \end{aligned} \quad (1)$$

and $\mathbf{v}(\mathbf{x})$ is estimated by minimizing Eq. 1 over a local 3D image patch $\Omega(\mathbf{x}, \mathbf{y})$, centered at \mathbf{x} . Note that v_t is not 1 since spatio-temporal orientation vectors will be computed. Using Lagrange multipliers, a corresponding error functional $e_{ls}(\mathbf{x})$ to minimize Eq. 1 using a least-squares error measure can be written as Eq. 2 where $W(\mathbf{x}, \mathbf{y})$ is a *spatially invariant* weighting function (e.g., Gaussian) that emphasizes the image gradients near the central pixel [8].

$$\begin{aligned} e_{ls}(\mathbf{x}) &= \int_{\Omega(\mathbf{x}, \mathbf{y})} (\nabla\mathbf{I}^T(\mathbf{y}) \mathbf{v}(\mathbf{x}))^2 W(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &\quad + \lambda \left(1 - \mathbf{v}(\mathbf{x})^T \mathbf{v}(\mathbf{x})\right) \end{aligned} \quad (2)$$

Assuming a constant $\mathbf{v}(\mathbf{x})$ within the neighborhood $\Omega(\mathbf{x}, \mathbf{y})$ and differentiating $e_{ls}(\mathbf{x})$ to find the minimum, leads to the standard eigenvalue problem (Eq. 3) for solving $\hat{\mathbf{v}}(\mathbf{x})$ the best estimate of $\mathbf{v}(\mathbf{x})$,

$$\mathbf{J}(\mathbf{x}, \mathbf{W}) \hat{\mathbf{v}}(\mathbf{x}) = \lambda \hat{\mathbf{v}}(\mathbf{x}) \quad (3)$$

The 3D structure tensor matrix $\mathbf{J}(\mathbf{x}, \mathbf{W})$ for the spatiotemporal volume centered at \mathbf{x} can be written in expanded matrix form, without the spatial filter $W(\mathbf{x}, \mathbf{y})$ and the positional terms shown for clarity, as Eq. 4.

$$\mathbf{J} = \begin{bmatrix} \int_{\Omega} \frac{\partial\mathbf{I}}{\partial x} \frac{\partial\mathbf{I}}{\partial x} d\mathbf{y} & \int_{\Omega} \frac{\partial\mathbf{I}}{\partial x} \frac{\partial\mathbf{I}}{\partial y} d\mathbf{y} & \int_{\Omega} \frac{\partial\mathbf{I}}{\partial x} \frac{\partial\mathbf{I}}{\partial t} d\mathbf{y} \\ \int_{\Omega} \frac{\partial\mathbf{I}}{\partial y} \frac{\partial\mathbf{I}}{\partial x} d\mathbf{y} & \int_{\Omega} \frac{\partial\mathbf{I}}{\partial y} \frac{\partial\mathbf{I}}{\partial y} d\mathbf{y} & \int_{\Omega} \frac{\partial\mathbf{I}}{\partial y} \frac{\partial\mathbf{I}}{\partial t} d\mathbf{y} \\ \int_{\Omega} \frac{\partial\mathbf{I}}{\partial t} \frac{\partial\mathbf{I}}{\partial x} d\mathbf{y} & \int_{\Omega} \frac{\partial\mathbf{I}}{\partial t} \frac{\partial\mathbf{I}}{\partial y} d\mathbf{y} & \int_{\Omega} \frac{\partial\mathbf{I}}{\partial t} \frac{\partial\mathbf{I}}{\partial t} d\mathbf{y} \end{bmatrix} \quad (4)$$

A typical approach in motion detection is to threshold $\text{trace}(\mathbf{J})$ (Eq. 5); but this results in ambiguities in distinguishing responses arising from stationary versus moving features (e.g., edges and junctions with and without motion), since $\text{trace}(\mathbf{J})$ incorporates total gradient change information but fails to capture the nature of these gradient changes (i.e. spatial only versus temporal).

$$\text{trace}(\mathbf{J}) = \int_{\Omega} \|\nabla I\|^2 dy \quad (5)$$

To resolve this ambiguity and to classify the video regions experiencing motion, the eigenvalues and the associated eigenvectors of \mathbf{J} are usually analyzed [10], [11]. However eigenvalue decompositions at every pixel is computationally expensive especially if real time performance is required.

B. Flux Tensors

In order to reliably detect only the moving structures *without* performing expensive eigenvalue decompositions, the concept of the *flux tensor* is proposed. Flux tensor is the temporal variations of the optical flow field within the local 3D spatiotemporal volume.

Computing the second derivative of Eq. 1 with respect to t , Eq. 6 is obtained where, $\mathbf{a}(\mathbf{x}) = [a_x, a_y, a_t]$ is the acceleration of the image brightness located at \mathbf{x} .

$$\frac{\partial}{\partial t} \left(\frac{d\mathbf{I}(\mathbf{x})}{dt} \right) = \frac{\partial^2 \mathbf{I}(\mathbf{x})}{\partial x \partial t} v_x + \frac{\partial^2 \mathbf{I}(\mathbf{x})}{\partial y \partial t} v_y + \frac{\partial^2 \mathbf{I}(\mathbf{x})}{\partial t^2} v_t + \frac{\partial \mathbf{I}(\mathbf{x})}{\partial x} a_x + \frac{\partial \mathbf{I}(\mathbf{x})}{\partial y} a_y + \frac{\partial \mathbf{I}(\mathbf{x})}{\partial t} a_t \quad (6)$$

which can be written in vector notation as,

$$\frac{\partial}{\partial t} (\nabla \mathbf{I}^T(x) \mathbf{v}(\mathbf{x})) = \frac{\partial \nabla \mathbf{I}^T(x)}{\partial t} \mathbf{v}(\mathbf{x}) + \nabla \mathbf{I}^T(x) \mathbf{a}(\mathbf{x}) \quad (7)$$

Using the same approach for deriving the classic 3D structure, minimizing Eq. 6 assuming a constant velocity model and subject to the normalization constraint $\|\mathbf{v}(\mathbf{x})\| = 1$ leads to Eq. 8,

$$e_{I_s}^F(\mathbf{x}) = \int_{\Omega(\mathbf{x}, \mathbf{y})} \left(\frac{\partial (\nabla \mathbf{I}^T(\mathbf{y}) \mathbf{v}(\mathbf{x}))}{\partial t} \right)^2 W(\mathbf{x}, \mathbf{y}) dy + \lambda \left(1 - \mathbf{v}(\mathbf{x})^T \mathbf{v}(\mathbf{x}) \right) \quad (8)$$

Assuming a constant velocity model in the neighborhood $\Omega(\mathbf{x}, \mathbf{y})$, results in the acceleration experienced by the brightness pattern in the neighborhood $\Omega(\mathbf{x}, \mathbf{y})$ to be zero at every pixel. As with its 3D structure tensor counterpart \mathbf{J} in Eq. 4, the 3D flux tensor \mathbf{J}_F using Eq. 8 can be written as

$$\mathbf{J}_F(\mathbf{x}, \mathbf{W}) = \int_{\Omega} W(\mathbf{x}, \mathbf{y}) \frac{\partial}{\partial t} \nabla \mathbf{I}(\mathbf{x}) \cdot \frac{\partial}{\partial t} \nabla \mathbf{I}^T(\mathbf{x}) dy \quad (9)$$

and in expanded matrix form as Eq. 10.

$$\mathbf{J}_F = \begin{bmatrix} \int_{\Omega} \left\{ \frac{\partial^2 \mathbf{I}}{\partial x \partial t} \right\}^2 dy & \int_{\Omega} \frac{\partial^2 \mathbf{I}}{\partial x \partial t} \frac{\partial^2 \mathbf{I}}{\partial y \partial t} dy & \int_{\Omega} \frac{\partial^2 \mathbf{I}}{\partial x \partial t} \frac{\partial^2 \mathbf{I}}{\partial t^2} dy \\ \int_{\Omega} \frac{\partial^2 \mathbf{I}}{\partial y \partial t} \frac{\partial^2 \mathbf{I}}{\partial x \partial t} dy & \int_{\Omega} \left\{ \frac{\partial^2 \mathbf{I}}{\partial y \partial t} \right\}^2 dy & \int_{\Omega} \frac{\partial^2 \mathbf{I}}{\partial y \partial t} \frac{\partial^2 \mathbf{I}}{\partial t^2} dy \\ \int_{\Omega} \frac{\partial^2 \mathbf{I}}{\partial t^2} \frac{\partial^2 \mathbf{I}}{\partial x \partial t} dy & \int_{\Omega} \frac{\partial^2 \mathbf{I}}{\partial t^2} \frac{\partial^2 \mathbf{I}}{\partial y \partial t} dy & \int_{\Omega} \left\{ \frac{\partial^2 \mathbf{I}}{\partial t^2} \right\}^2 dy \end{bmatrix} \quad (10)$$

As seen from Eq. 10, the elements of the flux tensor incorporate information about temporal gradient changes which leads to efficient discrimination between stationary and moving image features. Thus the trace of the flux tensor matrix which can be compactly written and computed as,

$$\text{trace}(\mathbf{J}_F) = \int_{\Omega} \left\| \frac{\partial}{\partial t} \nabla \mathbf{I} \right\|^2 dy \quad (11)$$

and can be directly used to classify moving and non-moving regions without the need for expensive eigenvalue decompositions. If motion vectors are needed then Eq. 8 can be minimized to get $\hat{\mathbf{v}}(\mathbf{x})$ using

$$\mathbf{J}_F(\mathbf{x}, \mathbf{W}) \hat{\mathbf{v}}(\mathbf{x}) = \lambda \hat{\mathbf{v}}(\mathbf{x}) \quad (12)$$

In this approach the eigenvectors need to be calculated at just moving feature points.

C. Flux Tensor Implementation

To detect motion blobs, only the trace of flux tensor $\text{trace}(\mathbf{J}_F) = \int_{\Omega(\mathbf{y})} \left\| \frac{\partial}{\partial t} \nabla \mathbf{I} \right\|^2 dy$ needs to be computed. That requires computation of I_{xt} , I_{yt} and I_{tt} and the integration of squares of I_{xt} , I_{yt} , I_{tt} over the area $\Omega(\mathbf{y})$. The following notation is adopted for simplicity:

$$I_{xt} = \frac{\partial^2 I}{\partial x \partial t}, \quad I_{yt} = \frac{\partial^2 I}{\partial y \partial t}, \quad I_{tt} = \frac{\partial^2 I}{\partial t^2} \quad (13)$$

The calculation of the derivatives is implemented as convolutions with a filter kernel. By using separable filters, the convolutions are decomposed into a cascade of 1D convolutions. For numerical stability as well as noise reduction, a smoothing filter is applied to the dimensions that are not convolved with a derivative filter e.g. calculation of I_{xt} requires smoothing in y-direction, and calculation of I_{yt} requires smoothing in x-direction. I_{tt} is the second derivative in temporal direction; the smoothing is applied in both spatial directions. As smoothing and derivative filters, optimized filter sets presented by Scharr et. al. in [12], [13] are used.

The integration is also implemented as an averaging filter decomposed into three 1D filters. As a result, calculation of trace at each pixel location requires three 1D convolutions for derivatives and three 1D convolutions for averages in the corresponding spatio-temporal cubes.

A brute-force implementation where spatial and temporal filters are applied for each pixel separately within a spatio-temporal neighborhood would be computationally very expensive since it would have to recalculate the convolutions for neighboring pixels. For an efficient implementation, the spatial (x and y) convolutions are separated from the temporal convolutions, and the 1D convolutions are applied to the whole frames one at a time. This minimizes the redundancy of computations and allows reuse of intermediate results.

The spatial convolutions required to calculate I_{xt} , I_{yt} and I_{tt} are I_{xs} , I_{sy} and I_{ss} where s represents the smoothing filter. Each frame of the input sequence is first convolved with two 1D filters, either a derivative filter in one direction and a smoothing filter in the other

direction, or a smoothing filter in both directions. These intermediate results are stored as frames to be used in temporal convolutions, and pointers to these frames are stored in a First In First Out (FIFO) buffer of size $n_{FIFO} = n_{Dt} + n_{At} - 1$ where n_{Dt} is the length of the temporal derivative filter and n_{At} is the length of the temporal averaging filter. For each input frame, three frames I_{xs} , I_{sy} and I_{ss} are calculated and stored. Once n_{Dt} frames are processed and stored, FIFO has enough frames for the calculation of the temporal derivatives I_{xt} , I_{yt} and I_{tt} . Since averaging is distributive over addition, $I_{xt}^2 + I_{yt}^2 + I_{tt}^2$ is computed first and spatial averaging is applied to this result and stored in the FIFO structure to be used in the temporal part of averaging. Once flux tensor trace of n_{At} frames are computed, temporal averaging is applied. Motion mask FG_M is obtained by thresholding and post-processing averaged flux tensor trace. Post-processing include morphological operations to join fragmented objects and to fill holes.

III. MOTION CONSTRAINED OBJECT SEGMENTATION USING GEODESIC ACTIVE CONTOURS

Motion blob detection produces a coarse motion mask FG_M as described in Section II-B. Each pixel of an infrared image frame $I_{IR}(x, t)$ is classified as moving or stationary by thresholding trace of the corresponding flux tensor matrix $\text{trace}(\mathbf{J}_F)$ and a motion blob mask $FG_M(t)$ is obtained.

Two problems with motion blob detection are: (1) holes: motion detection produces holes inside slow moving homogeneous objects, because of the aperture problem. (2) inaccurate object boundaries: motion blobs are larger than the corresponding moving objects, because these regions actually correspond to the union of the moving object locations in the temporal window, rather than the region occupied in the current frame. Beside inaccurate object boundaries this may lead to merging of neighboring object masks and consequently to false trajectory merges and splits at the tracking stage.

Motion constrained object segmentation module refines the coarse $FG_M(t)$ obtained through flux tensors by using fusion of multi-spectral image information and motion information, in a level set based geodesic active contours framework. Using mathematical morphology, holes in the motion blobs (disjoint foreground regions in FG_M) are filled. Geodesic active contours are started from the motion blob boundaries and evolved toward moving object boundaries defined by fusion of visible and infrared edges.

Contour evolution results in tighter object contours and separates most of the merged neighboring objects. Obtained object masks are further refined using shape information. Since the geodesic active contour segmentation relies on edges between background and foreground rather than the color or intensity differences such as in [14], the method is more stable and robust across very different appearances, non-homogeneous backgrounds and foregrounds. Starting the active contour evolution from

the motion segmentation results prevents early stopping of the contour on local non-foreground edges.

The motion constrained object segmentation process is summarized in Algorithm 1, level set based geodesic active contours, computation of edge indicator functions, shape-based segmentation refinement process are elaborated in the Sections III-A, IV, and V respectively.

Algorithm 1 Object Segmentation Algorithm

Input : Visible image sequence $I_{RGB}(x, t)$, infrared image sequence $I_{IR}(x, t)$, foreground mask sequence $FG_M(x, t)$ with $N_M(t)$ regions
Output : Refined foreground (binary) mask sequence $FG_R(x, t)$ with $N_R(t)$ regions

- 1: **for** each time t **do**
- 2: Compute edge indicator functions $g_{IR}(x, t)$ and $g_{RGB}(x, t)$ from infrared $I_{IR}(x, t)$ and visible $I_{RGB}(x, t)$ images.
- 3: Fuse $g_{IR}(x, t)$ and $g_{RGB}(x, t)$ into a single edge indicator function $g_F(x, t)$.
- 4: Initialize refined mask, $FG_R(t) \leftarrow 0$
- 5: Identify disjoint regions $R_i(t)$ in $FG_M(t)$ using connected component analysis.
- 6: **for** each region $R_i(t)$ $\{i = 1, 2, \dots, N_M(t)\}$ in $FG_M(t)$ **do**
- 7: Fill holes in $R_i(t)$ using morphological operations.
- 8: Initialize geodesic active contour level sets $C_i(t)$ using contour of $R_i(t)$.
- 9: Evolve $C_i(t)$ using $g_F(t)$ as edge stopping function.
- 10: Check stopping/convergence condition to subpartition $R_i(t) = \{R_{i,0}(t), R_{i,1}(t), \dots, R_{i,N_{R_i}(t)}(t)\}$ into $N_{R_i}(t) \geq 1$ foreground regions and one background region $R_{i,0}(t)$.
- 11: Refine mask FG_R using foreground partitions as $FG_R = FG_R \cup R_{i,j}; j = 1 : N_{R_i}(t)$
- 12: **end for**// $N_M(t)$ regions
- 13: **end for**// T frames

A. Level Set Based Geodesic Active Contours

Active contours evolve/deform a curve \mathcal{C} , subject to constraints from a given image. Active contours are classified as parametric [15], [16] or geometric [17]–[19] according to their representation. Parametric active contours (i.e. classical snakes) are represented explicitly as parametrized curves in a Lagrangian formulation, geometric active contours are implicitly represented as level sets [20] and evolve according to an Eulerian formulation [21]. Main advantages of level set based active contours over parametric active contours are computational simplicity and topological flexibility.

In this module level set based geodesic active contours [19] are used to refine motion blobs obtained using flux tensor method (Section II-B). These contours are topologically flexible and are effectively tuned to edge/contour information. Topological flexibility is critical in order to recover individual objects, because during the coarse motion segmentation, neighboring objects may have been merged into a single motion blob.

In level set based active contour methods the curve \mathcal{C} is represented implicitly via a Lipschitz function ϕ by $\mathcal{C} = \{(x, y) | \phi(x, y) = 0\}$, and the evolution of the curve is given by the zero-level curve of the function $\phi(t, x, y)$. Evolving \mathcal{C} in a normal direction with speed F amounts

to solving the differential equation [14],

$$\frac{\partial \phi}{\partial t} = |\nabla \phi| F; \quad \phi(0, x, y) = \phi_0(x, y) \quad (14)$$

In [17]–[19] geodesic active contour evolution is formulated as Eq. 15

$$\frac{\partial \phi}{\partial t} = g_F(\mathbf{I})(c + \mathcal{K}(\phi))|\nabla \phi| + \nabla \phi \cdot \nabla g_F(\mathbf{I}) \quad (15)$$

where $g_F(\mathbf{I})$ is the fused edge stopping function (Eq. 27), c is a constant, and \mathcal{K} is the curvature term,

$$\mathcal{K} = \text{div} \left(\frac{\nabla \phi}{|\nabla \phi|} \right) = \frac{\phi_{xx}\phi_y^2 - 2\phi_x\phi_y\phi_{xy} + \phi_{yy}\phi_x^2}{(\phi_x^2 + \phi_y^2)^{\frac{3}{2}}} \quad (16)$$

The force $(c + \mathcal{K})$ acts as the internal force in the classical energy based snake model. The constant velocity c pushes the curve inwards or outwards depending on its sign (inwards in our case). The regularization term \mathcal{K} ensures boundary smoothness. The external image dependent force $g_F(\mathbf{I})$ (Section IV) is the fused edge indicator function and is used to stop the curve evolution at visible or infrared object boundaries. The term $\nabla g_F \cdot \nabla \phi$ introduced in [19] is used to increase the basin of attraction for evolving the curve to the boundaries of the objects and to pull back the contour if it passes the boundary [21].

The level set function ϕ is initialized with the signed distance function of the motion blob contours (FG_M) and evolved using the geodesic active contour speed function Eq. 15 with edge stopping function $g_F(\mathbf{I})$ which fuses information from visible and infrared imageries. Formulation of edge stopping function and fusion of information from visible and infrared imageries are elaborated in the next section.

IV. VIDEO SENSOR FUSION COMBINING VISIBLE AND INFRARED EDGE INFORMATION

Contour feature or edge indicator functions are used to guide and stop the evolution of the geodesic active contour when it arrives at the object boundaries. The edge indicator function is a decreasing function of the image gradient that rapidly goes to zero along edges and is higher elsewhere.

The magnitude of the gradient of the infrared image is used to construct an edge indicator function g_{IR} as shown below where $G_\sigma(x, y) * \mathbf{I}_{IR}(x, y)$ is the infrared image smoothed with a Gaussian filter,

$$g_{IR}(\mathbf{I}_{IR}) = \exp(-|\nabla G_\sigma(x, y) * \mathbf{I}_{IR}(x, y)|) \quad (17)$$

Although the spatial gradient for single channel images lead to well defined edge operators, edge detection in multi-channel images (i.e. color edge strength) is not straight forward to generalize since gradients in different channels can have inconsistent orientations. In [22], Ruzon and Tomasi classify color edge detection algorithms into three categories based on when the individual channel responses are fused: (1) output fusion methods, (2) multi-dimensional gradient methods, and (3)

vector methods. In output fusion methods gray-scale edge detection is carried out in each channel independently then combined using methods such as weighted sum. Multi-dimensional gradient methods estimate a single estimate of the orientation and strength of an edge at a point. In vector methods, color information is treated as a vector through the whole process of edge detection such as in the case of edge detection based on vector order statistics [23], [24].

1) Tensor-based Color Edge Indicator Functions:

Color edge indicator functions in this work are based on multi-dimensional gradient methods. The main issue in multi-dimensional gradient methods is the combination of the individual channel gradients into a final multi-dimensional gradient. Simple methods use operations such as sum, weighted sum, max, min etc. to produce the final multi-dimensional/color gradients. But the summation of the individual channel gradients discard the correlation between the channels and may result in cancellation effects [25]. Pioneering work on how to combine the gradients of each channel is done by DiZenzo [26] who considered the multi-channel image as a vector field and computed the tensor gradient. Tensor based methods have since been used in various color feature detection algorithms. And many paper study and extend scale and affine invariance [27], [28] or photometric invariance [29] [30] of these features. More information on color features and multi-channel gradients can be found in [25] [22] [31]–[33]. Two types of tensors particularly important for multi-dimensional/color gradients are 2D color structure tensor and Beltrami color metric tensor explored below.

2D Color Structure Tensor:

The 2D color structure tensor is defined in Eq.18. Many color feature detectors [34]–[36] can be related to the eigenvalues of the color structure tensor matrix \mathbf{J}_C (Eq. 19). Since these eigenvalues are correlated with the local image properties of edgeness and cornerness i.e. $\lambda_1 \gg 0$, $\lambda_2 \approx 0$ and $\lambda_1 \approx \lambda_2 \gg 0$ respectively. Some local descriptors that can be obtained from the two eigenvalues Eq. 19 derived from the 2D color structure tensor are summarized in Table I [25].

$$\mathbf{J}_C = \begin{bmatrix} \sum_{i=R,G,B} \left(\frac{\partial \mathbf{I}_i}{\partial x} \right)^2 & \sum_{i=R,G,B} \frac{\partial \mathbf{I}_i}{\partial x} \frac{\partial \mathbf{I}_i}{\partial y} \\ \sum_{i=R,G,B} \frac{\partial \mathbf{I}_i}{\partial x} \frac{\partial \mathbf{I}_i}{\partial y} & \sum_{i=R,G,B} \left(\frac{\partial \mathbf{I}_i}{\partial y} \right)^2 \end{bmatrix} \quad (18)$$

$$\lambda_{1,2} = \frac{1}{2}(J_C(1,1) + J_C(2,2)) \pm \sqrt{(J_C(1,1) - J_C(2,2))^2 + (2J_C(1,2))^2} \quad (19)$$

$\lambda_1 + \lambda_2$	total local derivative energy
λ_1	derivative energy in the most prominent direction
$\lambda_1 - \lambda_2$	line energy
λ_2	amount of energy orthogonal to the prominent local orientation

TABLE I.: Local descriptors that can be obtained from the two eigenvalues λ_1, λ_2 derived from the 2D color structure tensor [25].

Beltrami Color Metric Tensor:

The Beltrami color metric tensor operator for a 2D color image defines a metric on a two-dimensional manifold $\{x, y, R(x, y), G(x, y), B(x, y)\}$ in the five-dimensional spatial-spectral space $\{x, y, R, G, B\}$. The color metric tensor is defined in Eq. 20 and reformulated as a function of the 2D color structure tensor in Eq. 21 where \mathcal{I}_2 is the 2×2 identity matrix and \mathbf{J}_C is the 2D color structure tensor. The magnitude of \mathcal{E} can be considered as a generalization of the gradient magnitude, and $\det(\mathcal{E})$ can be taken as the edge indicator function [37]–[39].

$$\mathcal{E} = \begin{bmatrix} 1 + \sum_{i=R,G,B} \left(\frac{\partial \mathbf{I}_i}{\partial x} \right)^2 & \sum_{i=R,G,B} \frac{\partial \mathbf{I}_i}{\partial x} \frac{\partial \mathbf{I}_i}{\partial y} \\ \sum_{i=R,G,B} \frac{\partial \mathbf{I}_i}{\partial x} \frac{\partial \mathbf{I}_i}{\partial y} & 1 + \sum_{i=R,G,B} \left(\frac{\partial \mathbf{I}_i}{\partial y} \right)^2 \end{bmatrix} \quad (20)$$

$$\mathcal{E} = \mathcal{I}_2 + \mathbf{J}_C \quad (21)$$

The Beltrami color edge stopping function can then be defined as,

$$g_{RGB}(\mathbf{I}_{RGB}) = \exp(-\text{abs}(\det(\mathcal{E}))) \quad (22)$$

$$\begin{aligned} \det(\mathcal{E}) &= \text{Beltrami}(\mathbf{I}_{RGB}) \\ &= 1 + \text{trace}(\mathbf{J}_C) + \det(\mathbf{J}_C) \\ &= 1 + (\lambda_1 + \lambda_2) + \lambda_1 \lambda_2 \end{aligned} \quad (23)$$

where λ_1 and λ_2 are the eigenvalues of \mathbf{J}_C .

Comparison of Color Feature Detectors:

Several common color (edge/corner) feature indicator functions were evaluated for comparison purposes, including Harris [34], Shi-Tomasi [35], Cumani [36] feature detectors and determinant of the Beltrami metric tensor [39].

The Harris operator (Eq. 24) [34] uses the parameter k to tune edge versus corner responses (i.e. $k \rightarrow 0$ responds primarily to corners).

$$\begin{aligned} \text{Harris}(\mathbf{I}_{RGB}) &= \det(\mathbf{J}_C) - k \text{trace}^2(\mathbf{J}_C) \\ &= \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2 \end{aligned} \quad (24)$$

The Shi-Tomasi operator (Eq. 25) [35] responds strongly to corners and filters out most edges (since one of the eigenvalues is nearly zero along edges). This is not suitable for a geodesic active contour edge stopping function.

$$\text{ShiTomas}(\mathbf{I}_{RGB}) = \min(\lambda_1, \lambda_2) \quad (25)$$

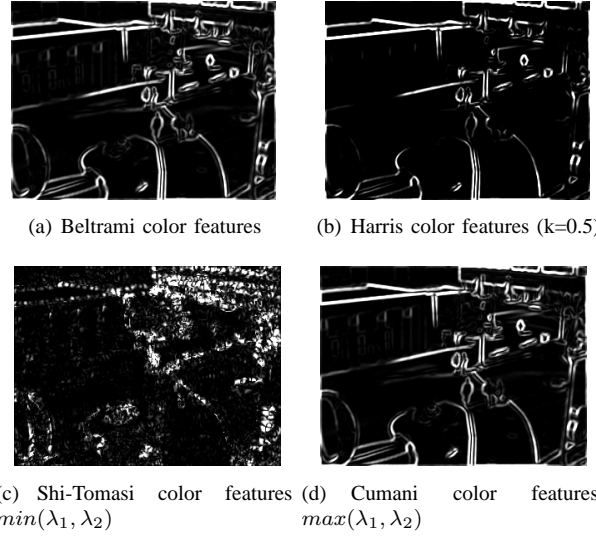


Figure 2: Color features for frame #1256 obtained using Beltrami, Harris, Shi-Tomasi, and Cumani operators.

The Cumani operator (Eq. 26) [36] responds nicely to both edges and corners.

$$\text{Cumani}(\mathbf{I}_{RGB}) = \max(\lambda_1, \lambda_2) \quad (26)$$

Although any robust accurate color edge response function can be used, the determinant of the Beltrami color tensor [39] has been selected for our application based on its robustness and speed. The best operators for the geodesic active contour edge stopping functions will respond to all salient contours in the image. In our experiments, as shown in Figure 2, the Beltrami color (edge) features was the most suitable function and is fast to compute. The Harris operator misses some salient contours around the pedestrians, the Shi-Tomasi operator responds primarily to corners and is not suitable as an edge stopping function, the Cumani operator produces a contour map that is nearly the same as the Beltrami color metric tensor map but is slightly more expensive to compute.

2) *Fusion of Infrared and Color Edge Indicator Functions:* Edge indicator functions obtained from infrared imagery and visible color imagery can be fused using various methods such as weighted sum or tensor-based fusion (i.e. infrared imagery can be considered as a fourth channel and the metric tensor in Eq. 18 can be defined in the six-dimensional spatial-spectral space $\{x, y, R, G, B, IR\}$). The fusion should satisfy two conditions: (1) fused edge stopping function $g_F(x, y)$ should respond to the strongest edge at location (x, y) in either IR or visible imagery, and (2) infrared imagery should have more weight in the final decision than any single channel of the visible imagery, since moving infrared edges are highly salient for tracking. In order to satisfy these conditions and not to miss any infrared edges, independent of the gradients in the visible channels, the *min* statistic Eq. 27 is used as the fusion operator. The fused edge stopping function g_F

is defined as the minimum of the two *normalized* (0, 1) edge stopping functions $g_{IR}(x, y)$, and $g_{RGB}(x, y)$,

$$g_F(\mathbf{IR}, \mathbf{RGB}, x, y) = \min\{g_{IR}(x, y), g_{RGB}(x, y)\} \quad (27)$$

This fusion method ensures that the curve evolution stops where there is an edge in the visible imagery or in the infrared imagery. *min* fusion operator handles cases where the visible RGB appearance of the moving object is similar to the background but there is a distinct infrared signature, and when the background and foreground have similar infrared signatures but distinct visible appearances.

V. SHAPE-BASED SEGMENTATION REFINEMENT

Evolving geodesic active contours from motion blob boundaries toward object boundaries defined by fusion of visible and infrared edges results in tighter object contours and separates most of the merged neighboring objects. But non-foreground edges such as background edges or shadow/illumination boundaries may result in early stopping of the active contours. Or close neighboring objects may lack non-edge spacing between them where the active contour can move. In such cases geodesic active contours result in under-segmentation of object clusters. Shape-based segmentation refinement breaks clusters in the refined foreground mask FG_R based on shape information. Shape-based refinement relies on an "ellipse like" shape model for people. Regions where people join are narrower compared to the size of the individual people that form the group/cluster.

The process is illustrated in Figure 3, summarized in Algorithm 2 and elaborated in the following paragraphs.

Algorithm 2 Shape-based Segmentation Refinement Algorithm

Input : $R \in FG_R(x, t)$ // Disjoint region in refined foreground mask
Output : R_i 's ($R = \cup_{i=1:n} R_i$, $n \geq 1$) // Sub-regions in R

```

1:  $R_E = R \ominus se_1$  // erode the region
2:  $L = Label(R_E)$  // apply connected component labeling
3:  $M_i \leftarrow$  Disjoint Regions in  $R_E$ ,  $R_E = \cup M_i$  // disjoint regions in  $R_E$  will be used as markers
4: if if ( $\#M_i \leq 1$ ) then
5:   Return R // keep R intact
6: else
7:   for each  $M_i \in R_E$  do
8:      $P_i = (M_i \oplus se_2) \cap FG_R$ ,  $se_2 > se_1$  // conditional dilation of the markers/parts
9:   end for
10:   $R_{rec} = \sum P_i$ 
11:   $Separators \leftarrow 0$ 
12:   $Separators(R_{rec} > 1) \leftarrow 1$  // overlapped regions are locations where the object should be splitted
13:   $Split \leftarrow ValidateSplit(R, Separators)$ 
14:  if if ( $Split == 0$ ) then
15:    Return R // keep R intact
16:  else
17:     $R(Separators == 1) \leftarrow 0$ 
18:     $R_i \leftarrow$  Disjoint Regions in R
19:    Return  $R_i$ 's
20:  end if
21: end if

```

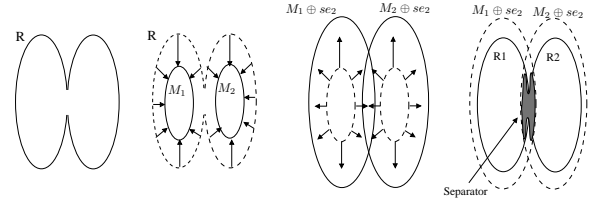


Figure 3: Shape-based segmentation refinement. R: original region/cluster, $M_{1,2}$: marker regions obtained by eroding R, $M_{1,2} \oplus se_2$: dilated markers, $R_{1,2}$: final segmentation.

Each disjoint region R of the refined mask FG_R is eroded and a set of markers (disjoint regions) is obtained. If R is a compact ellipse like region the erosion results in a single marker, no further processing is done and R is returned intact. If R contains narrow regions possibly corresponding to locations where people join, erosion removes those regions and the process results in $n > 1$ marker regions. R is partitioned using a process similar to morphological reconstruction. An n-layer image is allocated and each marker M_i is assigned to a layer. Similar to morphological reconstruction, successive conditional dilations (Eq. 28) are applied to the marker regions M_i 's (dilation is constrained to lie within the initial mask FG_R). While in reconstruction the process is stopped when the dilations cease to change, here the process is stopped when the dilated markers overlap with each other and the overlap regions are identified as partition separators.

$$P_i = (M_i \oplus se_2) \cap FG_R \quad (28)$$

Unlike morphological opening this reconstruction operation separates the objects but does not deteriorate their shapes. To avoid false fragmentation, validity of the partitioning is checked using an heuristic based on maximum, minimum, and relative sizes of the obtained partitions.

VI. MULTI-OBJECT TRACKING USING OBJECT CORRESPONDENCE GRAPHS

Moving object tracking is a fundamental step in the analysis of object behaviors [40]. Persistence in tracking is of particular importance for long term behavior analysis. Moving objects can be tracked implicitly i.e. through contour evolution [41] or explicitly through correspondence analysis. Active contour-based tracking methods track objects by evolving their contours dynamically in successive frames [40]. They provide an effective description of the tracked objects, but are highly sensitive to initialization and unless an additional re-initialization step such (e.g. feature-based re-initialization in [42]) is included to the process, they can not handle large displacements. Therefore, while active contours are used in the object segmentation stage (Section III), the proposed system adopts an explicit region correspondence based tracking algorithm that is able to handle larger object displacements.

The tracking module outlined in Algorithm 3 is an extension of our previous works in [43]–[45]. Object-to-object matching (correspondence) is performed using a

Algorithm 3 Tracking Algorithm

Input : Image sequence $\mathbf{I}(x, t)$, and refined foreground mask sequence $\text{FG}_R(x, t)$
Output : Trajectories and Temporal Object Statistics

- 1: **for** each frame $\mathbf{I}(x, t)$ at time t **do**
- 2: Use the refined foreground mask, $\text{FG}_R(t)$ from the motion constrained object segmentation module.
- 3: Partition $\text{FG}_R(t)$ into disjoint regions using connected component analysis $\text{FG}_R(t) = \{R_1(t), R_2(t), \dots, R_{N_R}(t)\}$ that ideally correspond to N_R individual moving objects.
- 4: **for** each region $R_i(t) \{i = 1, 2, \dots, N_R(t)\}$ in $\text{FG}_R(t)$ **do**
- 5: Extract blob centroid, area, bounding box, support map etc.
- 6: Arrange region information in an object correspondence graph \mathcal{O}_{GR} . Nodes in the graph represent objects $R_i(t)$, while edges represent object correspondences.
- 7: Search for potential object matches in consecutive frames using multi-stage overlap distance \mathcal{D}_{MOD} .
- 8: Update \mathcal{O}_{GR} by linking nodes that correspond to objects in frame $\mathbf{I}(x, t)$ with nodes of potential corresponding objects in frame $\mathbf{I}(x, t-1)$. Associate the confidence value of each match, $\mathcal{C}_{\mathcal{M}}(i, j)$ with each link.
- 9: **end for**
- 10: **end for**
- 11: Trace links in the object correspondence graph \mathcal{O}_{GR} to generate moving object trajectories.

multi-stage overlap distance \mathcal{D}_{MOD} , which consists of three distinct distance functions for three different ranges of object motion as described in [44].

Correspondence information is arranged in an acyclic directed graph \mathcal{O}_{GR} . Trajectory-Segments are formed by tracing the links of \mathcal{O}_{GR} and grouping "inner" nodes that have a single parent and a single child. For each Trajectory-Segment, parent and children segments are identified and a label is assigned. Segment labels encapsulate connectivity information and are assigned using a method similar to connected component labeling. Events such as appearance, disappearance, split and merge are also identified in this stage based on the number of parent and child segments.

Trajectory segment generation is followed by trajectory segment filtering and cluster trajectory analysis. Trajectory segment filtering prunes spurious trajectory segments based on features such as object size, trajectory length, displacement, and duration, or other spatio-temporal measures [43]. Factors such as under-segmentation, group interactions, partial or full occlusions result in temporary merging of individual object trajectories. Cluster trajectory analysis module resolves these merge-split events where n_p parent trajectory segments TPs, temporarily merge into a single trajectory segment TM, then split into n_c child trajectory segments TCs, and recovers individual object trajectories TSs. Currently only symmetric cases

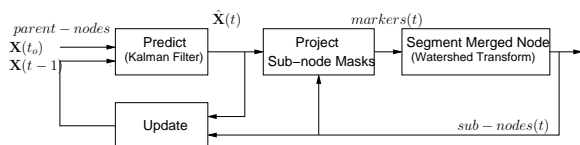


Figure 4: Cluster segmentation using Kalman filter and watershed segmentation.

Algorithm 4 Cluster Trajectory Analysis

Input : Merged trajectory segment TM, parent and child trajectory segments $\text{TP}(1:n_p)$, $\text{TC}(1:n_p)$.
Output : Individual trajectory segments $\text{TS}(1:n_p)$, updated parent and child trajectory segments $\text{TP}(1:n_p)$, $\text{TC}(1:n_p)$.

- 1: Initialize state matrix $\mathbf{X}(t_0)$ consisting of position and velocity informations for each individual trajectory segments $\text{TS}(1:n_p)$ using the parent segments' states, $\mathbf{X}(t_0) \leftarrow \text{TP}(1:n_p)$. states
- 2: **for** each node TM. node(t) of the merged segment **do**
- 3: Predict using Kalman filter [48] $\hat{\mathbf{X}}(t)$, the estimated state matrix of the trajectory segments $\text{TS}(1:n_p)$ corresponding to individual objects in the object cluster, from the previous states $\mathbf{X}(t-1)$.
- 4: Project masks of the sub-nodes $\text{TS}(1:n_p)$. node(t-1) from the previous frame, to the predicted positions on the current merged node TM. node(t), and use these projected masks as markers for the watershed segmentation.
- 5: Using watershed segmentation algorithm and the markers obtained through motion compensated projections of the sub-nodes from the previous frame, segment the merged node TM. node(t) corresponding to a cluster of objects into a set of sub-nodes corresponding to individual objects $\text{TS}(i)$. node(t), $i = 1 : n_p$.
- 6: Use refined positions of the sub-nodes obtained after watershed segmentation to update the corresponding states $\mathbf{X}(t)$.
- 7: **end for**
- 8: Update object correspondence graph \mathcal{O}_{GR} by including sub-node informations such as new support maps, centroids, areas etc.
- 9: Update individual trajectory segments's parent and children links ($\text{TS}(1:n_p)$. parents, $\text{TS}(1:n_p)$. children), parent segments' children links ($\text{TP}(1:n_p)$. children), and children segments' parent links ($\text{TC}(1:n_p)$. parents) by matching TSs to TPs and TCs.
- 10: Propagate parent segments' labels to the associated sub-segments, which are subsequently propagated to their children ($\text{TP}(1:n_p)$. label \rightarrow $\text{TS}(1:n_p)$. label \rightarrow $\text{TC}(1:n_p)$. label).

where $n_p = n_c$ are considered.

Most occlusion resolution methods rely heavily on appearance. But for far-view video, elaborate appearance-based models cannot be used since objects are small and not enough support is available for such models. Prediction and cluster segmentation is used to recover individual trajectories. Rather than predicting individual trajectories for the merged objects from the parent trajectories alone, at each time step, object clusters are segmented, new measurements are obtained, and object states are updated. This reduces error accumulation particularly for long lasting merges that become more frequent in persistent object tracking.

Segmentation of the object clusters into individual objects is done using a marker-controlled watershed segmentation algorithm applied to the object cluster masks [46], [47]. The use of markers prevents over-segmentation and enables incorporation of segmentation results from the previous frame. The cluster segmentation process is shown in Figure 4 and the overall cluster trajectory analysis process is summarized in Algorithm 4, where \mathbf{X} indicates state matrix that consists of a temporal sequence of position and velocity information for each individual object segment, and $\hat{\mathbf{X}}$ indicates estimated state matrix.

VII. RESULTS AND ANALYSIS

The proposed system is tested on thermal/color video sequence pairs from OTCBVS dataset collection [49].

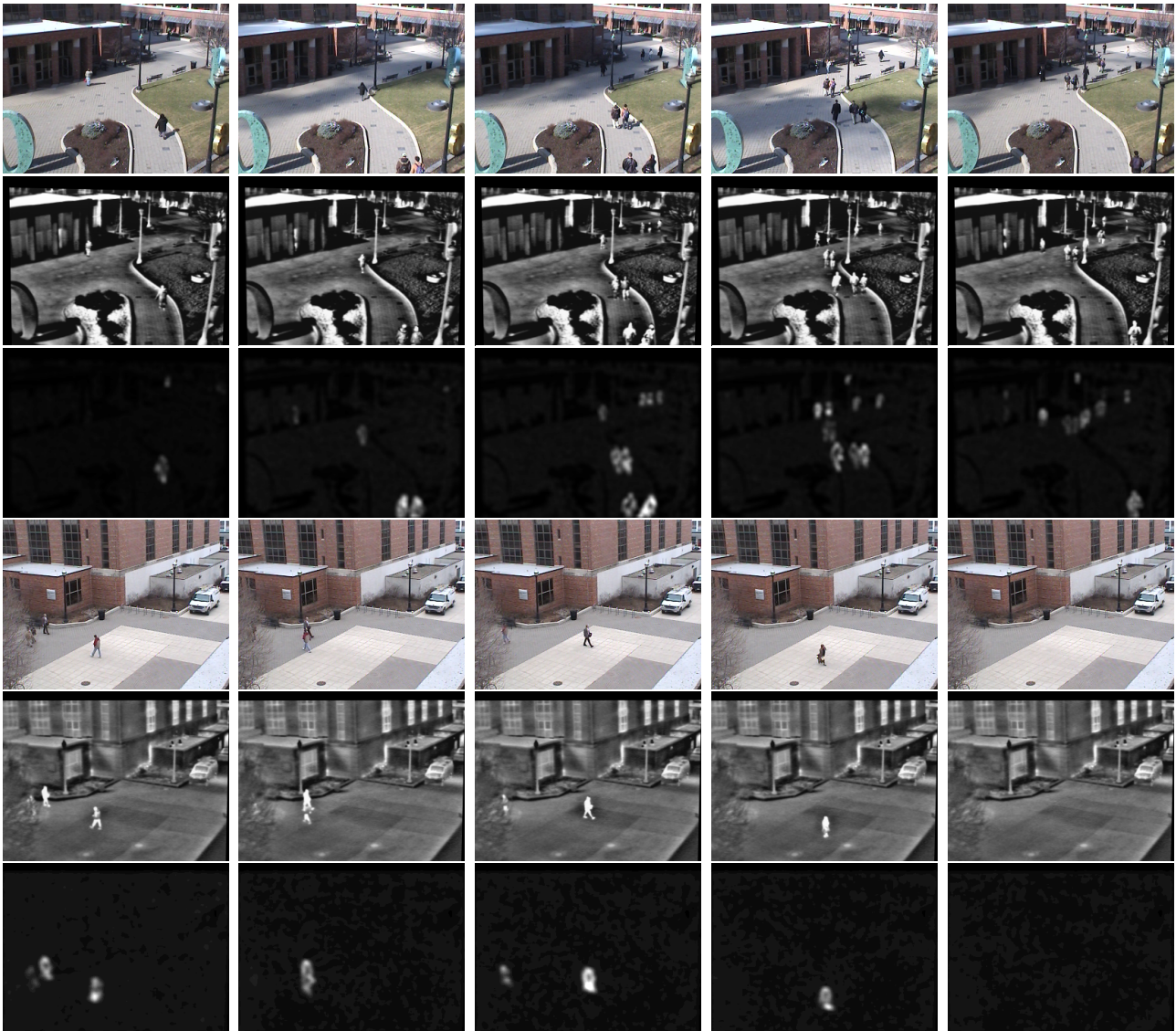


Figure 5: Flux tensor based motion detection results. Top three rows: visible image, IR image, flux tensor trace for OTCBVS benchmark sequence 3:1 frames #8,#432,#836,#1254,#1740. Flux tensor based motion detection results. Bottom three rows: visible image, IR image, flux tensor trace for OTCBVS benchmark sequence 3:4 frames #20,#124,#294,#1722,#3000. Averaging window size $n_{At} = 7$, derivative filter length $n_{Dt} = 5$. Flux response is scaled for improved print visibility.

Data consists of 8-bit grayscale bitmap thermal images, and 24-bit color bitmap images of 320 x 240 pixels. Images were sampled approximately at 30Hz. and registered using homography with manually-selected points. Thermal sequences were captured using a Raytheon PalmIR 250D sensor, color sequences were captured using a Sony TRV87 Handycam.

Figure 5 shows flux tensor trace for sample frames in OTCBVS sequences 3:1 and 3:4. While both IR image sequences particularly sequence 3:1 contain non-moving hot/bright regions such as part of the ground, roof tops, windows, lamp post etc., flux tensor trace successfully identifies only the moving pedestrians. Higher responses are produced for the object boundaries, and inside of the larger objects contain holes due to aperture problem.

Figure 6 shows different moving object detection results. MoG refers to the background estimation and subtraction method by mixture of Gaussians [50] [51]. Flux refers to the flux tensor method presented in Section II. The parameters for the mixture of Gaussians (MoG) method are selected as follows: number of distributions $K = 4$, distribution match threshold $T_{match} = 2.0$, background threshold $T = 70\%$, learning rate $\alpha = 0.02$. The parameters for the flux tensor method use a neighborhood size $W = 9$, and trace threshold $T = 4$. Visible imagery (Figure 6c) is very sensitive to moving shadows and illumination changes. Shadows (Figure 6c row 3) can alter object shapes and can result in false detections. Illumination changes due to cloud movements covers a large portion of the ground (Figure 6c row 4) which

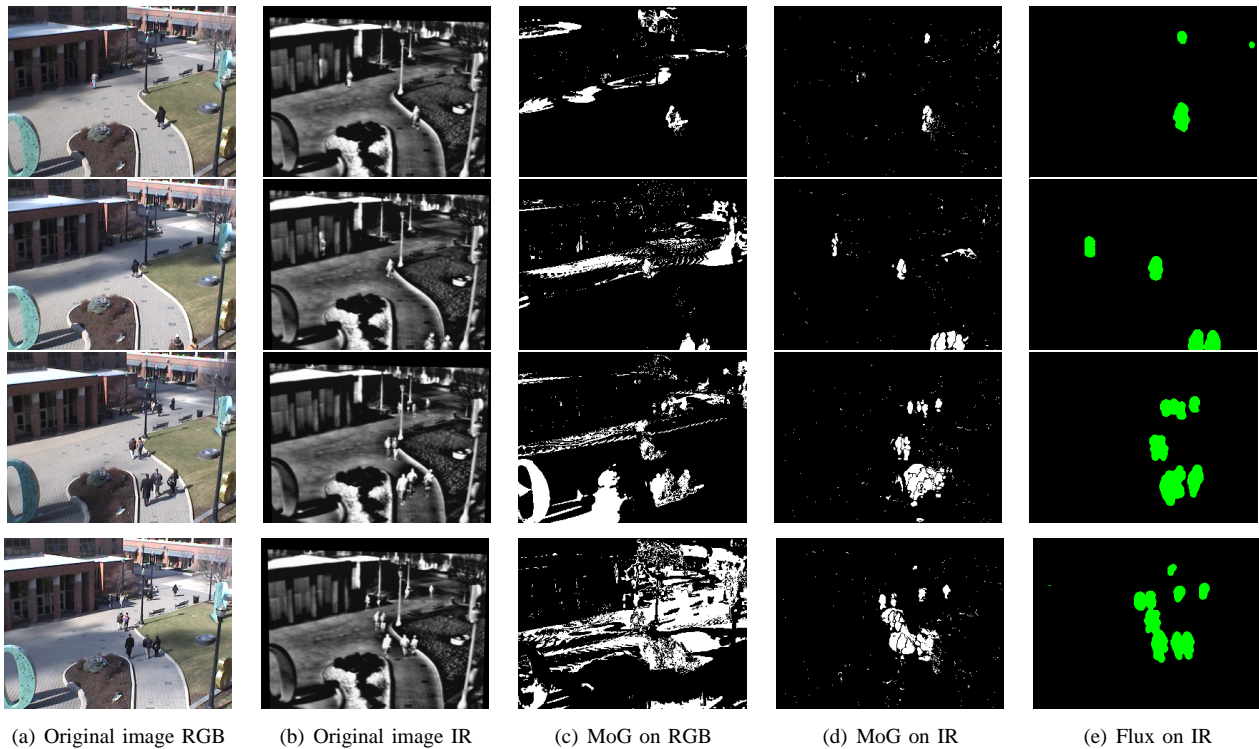


Figure 6: Moving object detection results for OTCBVS benchmark sequence 3:1. Top to bottom frames #120, #400,#1048,#1256.

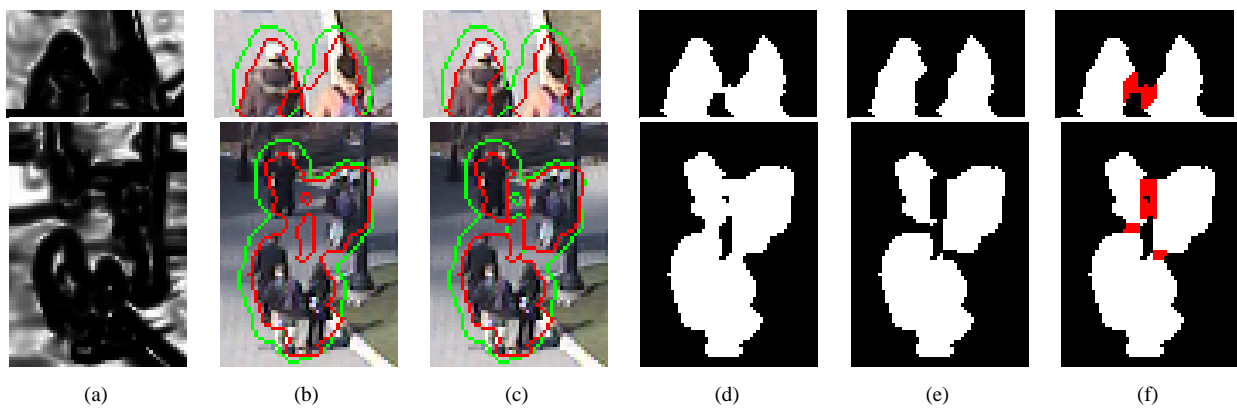


Figure 7: Shape-based segmentation refinement results for OTCBVS benchmark sequence 3:1 frames #408 (top) and #1528 (bottom). a) edge stopping function, b) contours before shape-based refinement, green: flux tensor results, red: contours refined using active contours, c) contours after shape-based refinement, green: flux tensor results, red: contours refined using active contour evolution + shape-based refinement, d) mask after active contour evolution, e) mask after active contour evolution + shape based refinement, f) partition separators.

results in many false moving object detections, making detection and tracking of pedestrians nearly impossible. As can be seen from Figures 6d,e infrared imagery is less sensitive to illumination related problems. But infrared imagery is more noisy compared to visible imagery and suffers from "halo" effects (Figure 6d). The flux tensor method (Figure 6e) produces less noisy and more compact foreground masks compared to pixel based background subtraction methods such as MoG (Figure 6d), since

it integrates temporal information from isotropic spatial neighborhoods.

Figure 7 illustrates two sample cases where the shape-based segmentation refinement splits the pedestrian clusters, under-segmented during active contour evolution. In the first case (top row), the two pedestrians walking side-by-side have no spacing between them where the geodesic active contour can move. In the second case (bottom row), the geodesic active contour stops early because of the

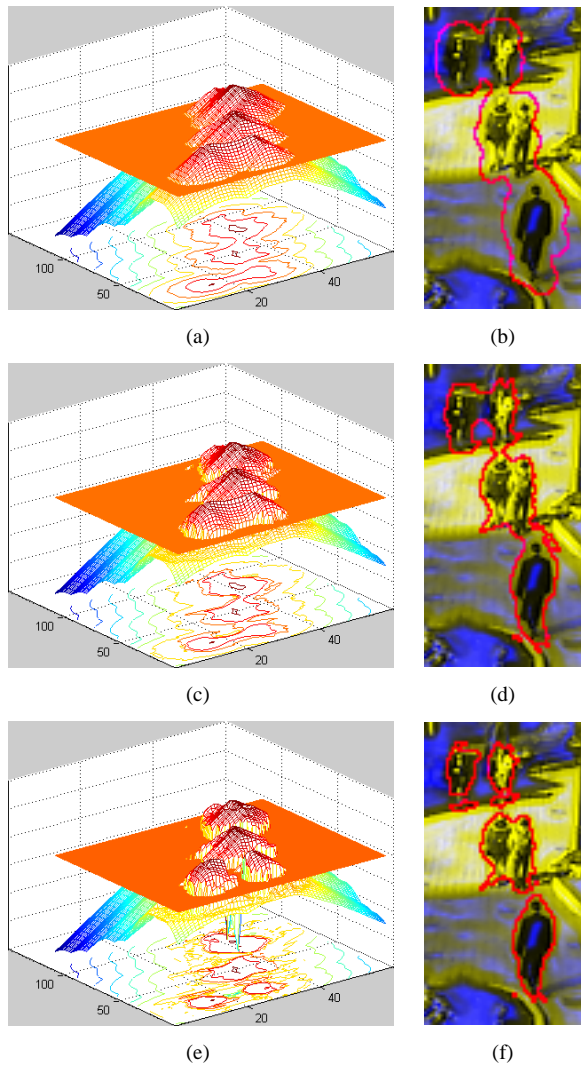
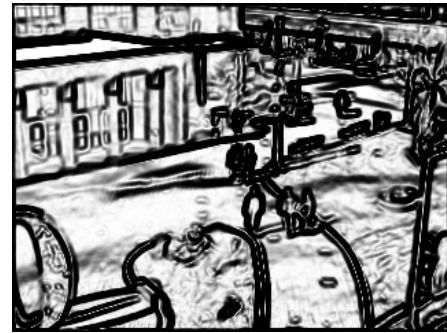


Figure 8: Level set based active contour evolution for OTCBVS benchmark sequence 3:1 frame #1256. Left: Level set function $\phi(x, t)$, horizontal plane is zero level. Right: Segmentation result superimposed on the original image. Red and Green channels: original image, Blue channel: fused edge stopping function $g_F(x, y)$, Red contour: level set contour ($\phi(x, y) = 0$).

background edges.

Figure 8 shows evolution of level set function $\phi(x, y)$ and the resulting contour during object segmentation. Level set evolution moves the contour from the motion blob boundaries inwards toward actual object boundaries defined by fusion of visible and infrared images. The process results in tighter object boundaries, and individual objects.

Figure 9 shows visible, IR and fused edge indicator functions used in the object segmentation process. Figure 10 illustrates effects of contour refinement and fusion of visible and infrared information. Level set based geodesic active contours refine object boundaries and segment object clusters into individual objects or smaller clusters



(a) g_{RGB}



(b) g_{IR}



(c) g_{Fusion}

Figure 9: Edge indicator functions for OTCBVS benchmark sequence 3:1 frame #1256.

which is critical for persistent object tracking. When used alone, both visible and infrared video result in total or partial loss of moving objects (i.e. top left person in Figure 10b due to low color contrast compared to background, parts of top right person and legs in Figure 10c due to lack of infrared edges). A low level fusion of the edge indicator functions shown in Figure 10d results in a more complete mask, compared to just combining visible and infrared foreground masks (i.e. legs of top right and bottom persons). Figure 11 illustrates the effects of contour refinement and merge resolution on object trajectories. Level set based geodesic active contours can separate clusters caused by under-segmentation (Figure 11a) but cannot segment individual objects during occlusions (Figure 11b). In those cases merge resolution recovers individual trajectories using prediction and previous object states (Figure 11 second row). In occlusion

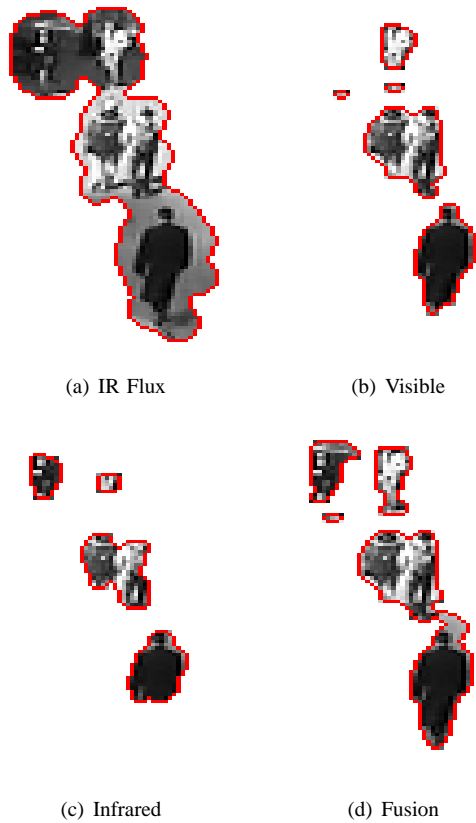


Figure 10: (a) Motion blob #2 in frame #1256 using IR flux tensors. Refinement of blob #2 using (b) only visible imagery, (c) only infrared imagery, (d) using fusion of both visible and IR imagery.

events no single parameter (i.e. color, size, shape etc.) can consistently resolve ambiguities in partitioning as evident in Figure 11b first row.

VIII. CONCLUSION

This paper presented a moving object detection and tracking system based on the fusion of infrared and visible imagery for persistent object tracking. Outdoor surveillance applications require robust systems due to wide area coverage, shadows, cloud movements and background activity. The proposed system fuses the information from both visible and infrared imagery within a geodesic active contour framework to achieve this robustness.

A new efficient motion detection algorithm referred to as the flux tensor is used to detect moving objects in infrared video without requiring background modeling or contour extraction. The flux tensor-based motion detector when applied to infrared video is more accurate than thresholding "hot-spots", and is insensitive to shadows as well as illumination changes in the visible channel. The novel flux tensor also produces less noisy and more spatially coherent results compared to classical pixel based background subtraction methods. The object segmentation algorithm uses level set-based geodesic active contour evolution that incorporates the fusion of visible color and

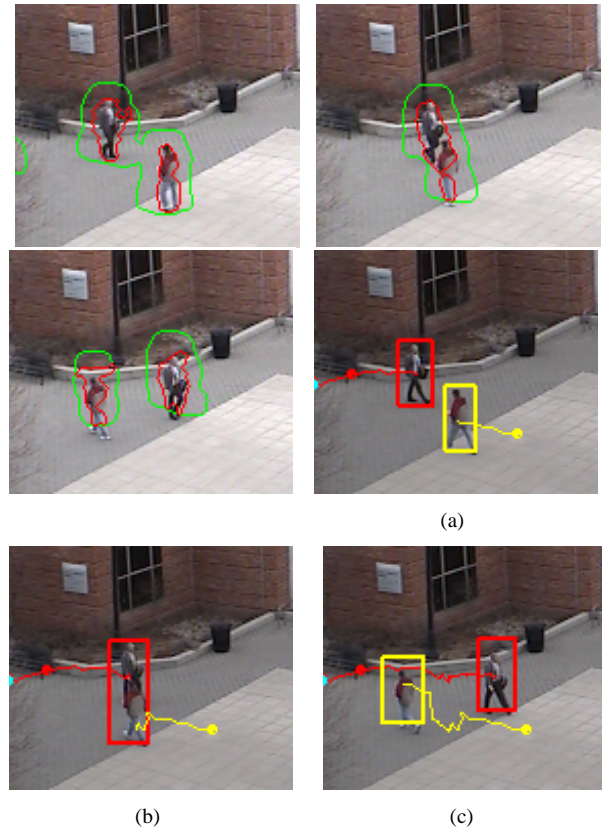


Figure 11: Merge resolution. Left to right: frames #41, #56, and #91 in OTCBVS benchmark sequence 3:4. Top row: motion constrained object extraction results. Flux tensor results are marked in green, refined contours are marked in red. Bottom row: object trajectories after merge resolution.

infrared edge information in a novel manner and refines the initial motion mask. Touching or overlapping objects are further refined during the segmentation process using an appropriate shape-based model.

Multiple object tracking using correspondence graphs is extended to handle groups of objects and occlusion events by Kalman filter-based cluster trajectory analysis and watershed segmentation. The proposed object tracking algorithm was successfully tested on several difficult outdoor multispectral videos from stationary sensors and is not confounded by shadows or illumination variations.

REFERENCES

- [1] J. Zolper, "Integrated microsystems: A revolution on five frontiers," in *Proc. of the 24th DARPA-Tech Conf.*, Anaheim, CA., Aug. 2005.
- [2] Z. Lemnios and J. Zolper, "Informatics: An opportunity for microelectronics innovation in the next decade," *IEEE Circuit & Devices*, vol. 22, no. 1, pp. 16–22, Jan. 2006.
- [3] G. Seetharaman, A. Lakhotia, and E. Blasch, "Unmanned vehicles come of age: The DARPA grand challenge," *Special Issue of IEEE Computer*, pp. 32–35, Dec. 2006.

- [4] W. Brown, R. Kaehr, and D. Chelette, "Finding and tracking targets: Long term challenges," *Air Force Research Technology Horizons*, vol. 5, no. 1, pp. 9–11, 2004.
- [5] F. Bunyak, K. Palaniappan, S. Nath, and G. Seetharaman, "Geodesic active contour based fusion of visible and infrared video for persistent object tracking," in *8th IEEE Workshop on Applications of Computer Vision (WACV 2007)*, Feb. 2007.
- [6] J. Davis and V. Sharma, "Background-subtraction in thermal imagery using contour saliency," *Int. Journal of Computer Vision*, vol. 71, no. 2, pp. 161–181, 2007.
- [7] S. Nath and K. Palaniappan, "Adaptive robust structure tensors for orientation estimation and image segmentation," in *LNCS-3804: Proc. ISVC'05*, Lake Tahoe, Nevada, Dec. 2005, pp. 445–453.
- [8] H.H. Nagel and A. Gehrke, "Spatiotemporally adaptive estimation and segmentation of OF-Fields," in *LNCS-1407: ECCV98*, Germany, June 1998, vol. 2, pp. 86–102.
- [9] B.P. Horn and B.G. Schunck, "Determining optical flow," *Artificial Intell.*, vol. 17, no. 1-3, pp. 185–203, Aug. 1981.
- [10] K. Palaniappan, H. Jiang, and T.I. Baskin, "Non-rigid motion estimation using the robust tensor method," in *IEEE Comp. Vision and Patt. Recog. Workshop on Articulated and Nonrigid Motion*, Washington, DC, June 2004.
- [11] J. Zhang, J. Gao, and W. Liu, "Image sequence segmentation using 3-D structure tensor and curve evolution," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 11, no. 5, pp. 629–641, May 2001.
- [12] H. Scharr, "Optimal filters for extended optical flow," in *LNCS: First Int. Workshop on Complex Motion*, Berlin, Germany, Oct. 2004, vol. 3417, pp. 66–74, Springer-Verlag.
- [13] H. Scharr, I. Stuke, C. Mota, and E. Barth, "Estimation of transparent motions with physical models for additional brightness variation," in *13th European Signal Processing Conference, EUSIPCO*, 2005.
- [14] T. Chan and L. Vese, "Active contours without edges," *IEEE Trans. Image Proc.*, vol. 10, no. 2, pp. 266–277, Feb. 2001.
- [15] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active contour models," *Journal International Journal of Computer Vision*, vol. 1, no. 4, pp. 321–331, January 1988.
- [16] L. D. Cohen, "On active contour models and balloons," *Computer Vision, Graphics, and Image Processing. Image Understanding*, vol. 53, no. 2, pp. 211–218, 1991.
- [17] V. Caselles, F. Catte, T. Coll, and F. Dibos, "A geometric model for active contours," *Numerische Mathematik*, vol. 66, pp. 1–31, 1993.
- [18] R. Malladi, J.A. Sethian, and B. Vemuri, "Shape modelling with front propagation: A level set approach," *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. 17, no. 2, pp. 158–174, 1995.
- [19] V. Caselles, R. Kimmel, and G. Sapiro, "Geodesic active contours," *Int. Journal of Computer Vision*, vol. 22, no. 1, pp. 61–79, 1997.
- [20] J.A. Sethian, *Level Set Methods and Fast Marching Methods: Evolving Interfaces in Computational Geometry, Fluid Mechanics, Computer Vision, and Materials Science*, Cambridge University Press, Cambridge, UK, 1999, ISBN 0-521-64557-3.
- [21] C. Xu, A. Yezzi, and J. Prince, "A summary of geometric level-set analogues for a general class of parametric active contour and surface models," in *Proc. IEEE Workshop on Variational and Level Set Methods in Computer Vision*, 2001, pp. 104–111.
- [22] M.A. Ruzon and C. Tomasi, "Edge, junction, and corner detection using color distributions," *IEEE Trans. Patt. Anal. and Mach. Intell.*, vol. 23, no. 11, pp. 1281–1295, 2001.
- [23] P. E. Trahanias and A. N. Venetsanopoulos, "Vector order statistics operators as color edge detectors," *IEEE Transactions on Systems, Man, and Cybernetics-Part B: Cybernetics*, vol. 26, no. 1, pp. 135–143, Feb. 1996.
- [24] Hai Tao and Thomas S. Huang, "Color image edge detection using cluster analysis," in *IEEE Int. Conf. on Image Processing (ICIP'97)*, 1997, pp. 834–836.
- [25] H. Stokman T. Gevers, J. Weijer, "Color feature detection," in *Color Image Processing: Methods and Applications*, R. Lukac and K.N. Plataniotis, Eds. CRC Press, 2006.
- [26] Silvano Di Zenzo, "A note on the gradient of a multi-image," *Computer Vision, Graphics, and Image Processing*, vol. 33, no. 1, pp. 116–125, 1986.
- [27] K. Mikolajczyk and C. Schmid, "Scale & affine invariant interest point detectors," *International Journal on Computer Vision*, vol. 60, no. 1, pp. 63–86, 2004.
- [28] Nicu Sebe, Theo Gevers, S. Dijkstra, and Joost van de Weijer, "Evaluation of intensity and color corner detectors for affine invariant salient regions," in *Beyond Patches Workshop Interaction at CVPR*, 2006.
- [29] J. van de Weijer, Th. Gevers, and J.M. Geusebroek, "Edge and corner detection by photometric quasi-invariants," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 4, April 2005.
- [30] J. van de Weijer, Th. Gevers, and A.W.M. Smeulders, "Robust photometric invariant features from the color tensor," *IEEE Trans. Image Proces.*, vol. 15, no. '1, 2006.
- [31] C. Schmid, R. Mohr, and C. Bauckhage, "Evaluation of interest point detectors," *International Journal of Computer Vision*, vol. 37, no. 2, pp. 151–172, 2000.
- [32] A. Koschan and M. Abidi, "Detection and classification of edges in color images," *Signal Processing Magazine, Special Issue on Color Image Processing*, vol. 22, no. 1, pp. 64–73, 2005.
- [33] S. K. Naik and C. A. Murthy, "Standardization of edge magnitude in color images," *IEEE Trans. on Image Processing*, vol. 15, no. 9, pp. 2588–2595, Sep. 2006.
- [34] C. Harris and M. Stephens, "A combined corner and edge detector," in *Proc. 4th Alvey Vision Conf.*, Manchester, 1988, vol. 15, pp. 147–151.
- [35] Jianbo Shi and Carlo Tomasi, "Good features to track," in *IEEE Conf. on Comp. Vis. and Patt. Recog.(CVPR)*, Seattle, June 1994.
- [36] A. Cumani, "Edge detection in multispectral images," *CVGIP: Graphical Models and Image Processing*, vol. 53, no. 1, 1991.
- [37] A. Brook, R. Kimmel, and N.A. Sochen, "Variational restoration and edge detection for color images," *Journal of Mathematical Imaging and Vision*, vol. 18, no. 3, pp. 247–268, May 2003.
- [38] N.A. Sochen, R. Kimmel, and R. Malladi, "A general framework for low-level vision," *IEEE Trans. Image Proc.*, vol. 7, no. 3, pp. 310–318, 1998.
- [39] R. Goldenberg, R. Kimmel, E. Rivlin, and M. Rudzsky, "Fast geodesic active contours," *IEEE Trans. Image Proc.*, vol. 10, no. 10, pp. 1467–1475, Oct 2001.
- [40] W. Hu, T.N. Tan, L. Wang, and S.J. Maybank, "A survey on visual surveillance of object motion and behaviors," *IEEE Trans. on Systems, Man, and Cybernetics - Part C*, vol. 34, no. 3, pp. 334–352, August 2004.
- [41] N. Paragios and R. Deriche, "Geodesic active regions and level set methods for motion estimation and tracking," *Computer Vision and Image Understanding*, vol. 97, no. 3, pp. 259–282, March 2005.

- [42] A. Knoll G. Panin, "Fully automatic real-time 3d object tracking using active contour and appearance models," *Journal of Multimedia* 2006, vol. 1, no. 7, pp. 62–70, 2006.
- [43] F. Bunyak and S. R. Subramanya, "Maintaining trajectories of salient objects for robust visual tracking," in *LNCS-3212: Proc. ICIAR'05*, Toronto, Sep. 2005, pp. 820–827.
- [44] F. Bunyak, K. Palaniappan, S.K. Nath, T.I. Baskin, and G. Dong, "Quantitive cell motility for *in vitro* wound healing using level set-based active contour tracking," in *Proc. 3rd IEEE Int. Symp. Biomed. Imaging (ISBI)*, pp. 1040–1043, Arlington, VA, April 2006.
- [45] S. K. Nath, F. Bunyak, and K. Palaniappan, "Robust tracking of migrating cells using four-color level set segmentation," in *LNCS-3212: Proc. ACIVS'05*, Antwerp, Belgium, Sep. 2006, vol. 3212.
- [46] L. Vincent and P. Soille, "Watersheds in digital spaces: an efficient algorithm based on immersion simulations," *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. 13, no. 6, pp. 583–598, 1991.
- [47] S. Beucher and F. Meyer, "The morphological approach to segmentation: the watershed transformation," in *Math. Morph. and its Applications to Image Proc.*, E.R. Dougherty, Ed., pp. 433–481. Marcel Dekker, NY, 1993.
- [48] Y. Bar-Shalom, X.R. Li, and T. Kirubarajan, *Estimation with Applications to Tracking and Navigation: Theory, Algorithms, and Software*, John Wiley & Sons, Inc., 2001.
- [49] J. Davis and V. Sharma, "Fusion-based background-subtraction using contour saliency," in *IEEE Int. Workshop on Object Tracking and Classification Beyond the Visible Spectrum*, San Diego, CA, June 2005.
- [50] C. Stauffer and E. Grimson, "Learning patterns of activity using real-time tracking," *IEEE Trans. Pattern Anal. and Machine Intel.*, vol. 22, no. 8, pp. 747–757, 2000.
- [51] X. Zhuang, Y. Huang, K. Palaniappan, and Y. Zhao, "Gaussian mixture density modeling, decomposition and applications," *IEEE Trans. Image Proc.*, vol. 5, no. 9, pp. 1293–1302, Sep. 1996.

Filiz Bunyak received her B.S. and M.S. degrees in control and computer engineering from Istanbul Technical University, Istanbul, Turkey and her Ph.D. degree in computer science from University of Missouri-Rolla, Rolla, MO, USA in 2005. In 2005 she joined University of Missouri-Columbia computer science department as post-doctoral researcher. Her research interests are image processing and computer vision, with emphasis on visual surveillance, video tracking, biomedical image processing, data fusion, segmentation, level set methods, and mathematical morphology.

Kannappan Palaniappan received the B.A.Sc. and M.A.Sc. degrees in systems design engineering from the University of Waterloo, Waterloo, ON, Canada, and the Ph.D. degree in electrical and computer engineering from the University of Illinois, Urbana-Champaign, in 1991. From 1991 to 1996, he was with the NASA Goddard Space Flight Center, working in the Laboratory for Atmospheres, where he coestablished the High-Performance Visualization and Analysis Laboratory. He developed the first massively parallel algorithm for satellite-based hurricane motion tracking using 64 K processors and invented the Interactive Image Spreadsheet system for visualization of extremely large image sequences and numerical model data over highperformance networks. Many visualization products created with colleagues at NASA have been widely used on television, magazines, museums, web sites, etc. With the University of Missouri, he helped establish the NSF vBNS highspeed research network, the NASA Center of Excellence in Remote Sensing, ICREST, and MCVL. He has been with UMI-ACS, University of Maryland, and worked in industry for Bell

Northern Research, Bell Canada, Preussen Elektra Germany, the Canadian Ministry of Environment, and Ontario Hydro. His research interests include satellite image analysis, biomedical imaging, video tracking, level set-based segmentation, nonrigid motion analysis, scientific visualization, and content-based image retrieval. Dr. Palaniappan received the highest teaching award given by the University of Missouri, the William T. Kemper Fellowship for Teaching Excellence in 2002, the Boeing Welliver Faculty Fellowship in 2004, the University Space Research Association Creativity and Innovation Science Award (1993), the NASA Outstanding Achievement Award (1993), the Natural Sciences and Engineering Research Council of Canada scholarship (1982/1988), and the NASA Public Service Medal (2001) for pioneering contributions to scientific visualization and analysis tools for understanding petabyte-sized archives of NASA datasets.

Sumit Nath received his M.Tech degree from the Indian Institute of Science, Bangalore, India in 1998 and his Phd from the University of Ottawa, Canada, 2004. He joined Dr. K.Palaniappan's group at the University of Missouri-Columbia as a Post-Doc researcher in 2004. Currently, he is a pursuing his post-doc research at the Electrical Computer and Systems Engineering Department at Rensselaer Polytechnic Institute, Troy. He is primarily interested in all aspects of computer vision, image processing and compression, with special emphasis on level sets, segmentation, optic flow estimation, Voronoi diagrams and Tomography.

Guna Seetharaman has been an associate professor of Electrical and Computer Engineering at Air Force Institute of Technology, since June 2003. He was with The Center for Advanced Computer Studies, University of Louisiana at Lafayette, until 2003, as an associate professor of computer engineering. He was a CNRS Research Visiting Professor at the Institute for Electronics Fundamentals, University of Paris XI, Orsay, on sabbatical and short-term fellowships.

Dr. Seetharaman earned his Ph.D. degree in Electrical and Computer Engineering, in 1988 from the University of Miami, Coral Gables, FL. He holds a M.Tech., in Electrical Engineering (1982) from Indian Institute of Technology, Chennai. He earned his B.E., Electronics and Communication Engineering, in 1980 from The University of Madras, India. He established and successfully ran the Computer Vision Laboratory, and Intelligent Robotics Laboratory (co-established) at The University of Louisiana at Lafayette. His research has been funded by NSF, ONR, DOE, AFOSR and The Board of Regents of Louisiana. His recent focus has been on integrated Microsystems for 3D imaging and displays; and high performance embedded computing algorithms for image processing systems. He also participated in the DARPA Grand Challenge as a charter member of Team CajunBot. He has published more than 120 articles in: Computer Vision, low-altitude aerial imagery, SIMD-Parallel Computing, VLSI-signal processing, 3D Displays, Nano-Technology, and 3D Video analysis.

He co-organized the DOE/ONR/NSF Sponsored Second International Workshop on Foundations of Decision and Information Fusion, in 1996 (Washington DC), and the IEEE Sixth International Workshop on Computer Architecture for Machine Perception, New Orleans, 2003. He guest edited IEEE COMPUTER special issue devoted to Unmanned Intelligent Autonomous Vehicles, Dec 2006. He also guest-edited a special issue of the EURASIP Journal on Embedded Computing in the topics of Intelligent Vehicles. He is an active member of the IEEE, and ACM. He is also a member of Tau Beta Pi, Eta Kappa Nu and Upsilon Pi Epsilon.