

# Online Real-time Crowd Behavior Detection in Video Sequences

Andrea Pennisi<sup>a</sup>, Domenico D. Bloisi<sup>a,\*</sup>, Luca Iocchi<sup>a</sup>

<sup>a</sup>*Department of Computer, Control, and Management Engineering  
Sapienza University of Rome  
via Ariosto 25, 00185, Rome (Italy)*

---

## Abstract

Automatically detecting events in a crowded scene is a challenging task in computer vision. A number of offline approaches have been proposed for solving the problem of crowd behavior detection, however the offline assumption limits their application in real video surveillance systems. In this paper, we propose an online and real-time method for detecting events in video sequences. The proposed approach is based on the combination of visual feature extraction and image segmentation and it works without the need of a training stage. A quantitative experimental evaluation carried out on publicly available video sequences, containing data from various crowd scenarios and different types of events, demonstrates the effectiveness of the approach.

*Keywords:* event detection, crowd analysis, image segmentation, intelligent surveillance

---

## 1. Introduction

Event detection in the field of automatic video surveillance has gained a growing interest [1]. The huge amount of data generated by existing surveillance systems requires the development of more intelligent systems with the final goal of avoiding information overload for the users [2]. In particular, in the context of a crowd image analysis problem, it is desirable to develop

---

\*Corresponding author

*Email address:* [bloisi@dis.uniroma1.it](mailto:bloisi@dis.uniroma1.it) (Domenico D. Bloisi)

*URL:* <http://www.dis.uniroma1.it/~bloisi/> (Domenico D. Bloisi)

online algorithms that reliably detect abnormal events in real-time. For example, the automatic analysis of crowded scenes can be used to avoid crowd related disasters and ensure public safety [3].

An anomaly can be defined as: “*something that deviates from what is standard, normal, or expected*”<sup>1</sup>. It means that *abnormal* events can be identified as irregular events with respect to usual normal ones. Thus, the abnormal detection becomes the identification of abnormal events given some sample normal events. Zhan et al. in [4] point out that conventional computer vision can be ineffective when dealing with the analysis of very crowded video sequences. Indeed, in a high-density situation the presence of severe occlusions consistently limits the performance of traditional methods for visual tracking [3]. Additional factors that can limit the effectiveness of existing approaches aiming at detecting abnormal events are:

- Off-line computation;
- Need of a training phase.

The off-line assumption can limit the application of the anomaly detection method in practice [5]. For instance, it is desirable to detect panic situations as soon as possible in order to avoid damage to people. Methods that relies on a training phase are limited by the possible lack of well-suited training data. Indeed, since it is not easy to find data about real emergency situations in crowded scenes, the resulting classifier could be suitable only for dealing with particular video sequences.

In this paper, we propose an online, real-time method for automatic anomaly detection in crowded scenes, that does not need any training stage. In particular, the main contributions of the proposed approach are:

1. The definition of two different metrics, namely *instant entropy* and *temporal occupancy variation*, to detect abnormal situations in crowded scenes;
2. A segmentation algorithm for images containing crowds.

Furthermore, we provide:

- A novel video sequence annotated with ground truth data, containing images of hundreds of runners at the start of a marathon, as an example of crowd video with locally steady optical flow.

---

<sup>1</sup>Definition from the Oxford Dictionary.

- Ground truth data for two well-know video sequences containing crowd scenes, namely PETS 2009 [6] and AGORASET [7].
- The source code and all the data used for the experimental evaluation at the following web page <http://www.dis.uniroma1.it/~pennisi/eventdetection.html>, thus allowing for reproducing the results described in this paper and to compare other similar approaches.

The reminder of the paper is organized as follows. Related work is analyzed in the next Section 2, while our method is presented in Section 3. Section 4 describes the qualitative and quantitative experimental results, providing also a comparison with other approaches in the literature. Conclusions are provided in Section 5.

## 2. Related Work

Techniques for crowd behavior analysis are usually grouped into two main categories [3, 8]: object-based and holistic approaches. In the object-based methods the analysis is carried out at an individual level. For example, it can be of interest to detect if a single person is trying to enter a restricted area or if an individual is moving against the dominant flow. On the other hand, holistic techniques treats the crowd as a single entity, trying to extract global information, such as the main flow of the crowd, instead of analyzing single trajectories.

In this section, we provide a different classification based on the nature of the methods used for detecting abnormal situations. In particular, existing approaches are grouped into:

- Statistical analysis;
- Background subtraction;
- Segmentation;
- Classification.

*Statistical analysis.* Methods in this category are based on the collection of particular features representing the flow of the crowd. For example, Zhang et al. in [9] describe a social attribute-aware force model for abnormal crowd pattern detection in video sequences. An unsupervised method is used to estimate the scene scale and a social disorder attribute and congestion attribute

are introduced to describe the realistic social behaviors by using statistical context feature. Through the semantic attribute-aware enhancement, they obtain an improved model on the basis of social force. Even if the method has good results, it is an off-line method.

Zhu and Saligrama in [1] propose a probabilistic framework that takes into consideration local spatio-temporal anomalies in order to characterize the observed scene by using optimal decision rules. If anomalies are local optimal decision, they are local as well, even if the behavior exhibits global spatial and temporal statistical dependencies. This helps to collapse the large ambient data dimension space in order to detect local anomalies. Consistent data-driven local empirical rules with provable performance can be derived with limited training data. The empirical rules are based on scores functions derived from local nearest neighbor distances. These rules aggregate statistics across spatio-temporal locations and scales, and produce a single composite score for video segments.

Chang et al. [10] describe a statistical framework able to recognize group-level activity in many scenarios, using a soft grouping metric and track-based motion analysis. The approach recognizes group interactions without making hard decisions about the underlying group structure. In particular, a path-based grouping scheme is used to understand if an individual belongs to a group. The method is bottom up and thus could be limited where the tracking is not reliable.

Mehran et al. [8] propose a method for localizing abnormal behaviors by using a Social Force model. A grid of particles is placed over the image for analyzing the space-time average of optical flow. The moving particles are treated as individuals and the social forces are estimated by using the social force model. The interaction forces are then mapped into the image plane to obtain Force Flow for every pixel in every frame. Spatio-temporal volumes of Force Flow are randomly selected for modeling the normal behavior of the crowd. Then, the normal and abnormal behaviors are classified by using an approach based on a bag of words. The regions of anomalies in the abnormal frames are localized using interaction forces.

Kratz and Nishino [11] describe a statistical framework for modeling the motion pattern behavior of extremely crowded scenes in order to detect unusual events. The authors model the dense activity of the crowd using a 3D Gaussian distribution of spatio-temporal gradients, capturing the local spatio-temporal motion patterns through a distribution-based Hidden Markov Model. The results demonstrate that the used approach is a suit-

able representation for analyzing crowded scenes, detecting unusual motion patterns in pedestrian behavior including movement against the normal flow of traffic.

*Background Subtraction.* Approaches that uses background subtraction are commonly based on the creation of a Gaussian Mixture Model (GMM) to extract foreground objects. For example, Fradi et al. in [12] propose a people counting approach that harness the advantage of incorporating an uniform motion model into GMM background subtraction to obtain high accurate foreground segmentation. The counting is based on foreground measurements, where a perspective normalization and a crowd measure-informed corner density are introduced with foreground pixel counts into a single feature. The approach demonstrates the benefits of integrating GMM with motion cue and normalizing the proposed feature as well. However, it is not adaptive to illumination conditions.

Srivastava et al. in [13] describe a method for crowd flow estimation by counting the number of persons passing through a designated region in a unit time. The method accumulates the total number of foreground pixels over a chosen time period that is directly proportional to the number of people passing to the defined area through a scaling factor. This factor depends on the local texture features that takes into account the level of occlusions.

Li et al. in [14] propose a foreground detection approach for crowd motion analysis called optical flow and background model (OFBM) that relies on Lucas-Kanade optical flow and Gaussian background model methods to eliminate the noise due to brightness changes and occlusions. This approach overcomes the shortages of optical flow and background subtract, but it is not computationally fast enough to be applied in real-time processing.

*Segmentation.* Methods in this category rely on the identification of the crowd flow by using a grid particles placed in the scene in order to detect the evolution of the people in the scene. Solmaz et al. in [15] propose a framework to identify multiple crowd behaviors through stability analysis for dynamical systems. A scene is overlaid by a grid of particles initializing a dynamical system defined by the optical flow. Time integration of the dynamical system provides particle trajectories that represent the motion in the scene; then, these trajectories are used to locate regions of interest in the scene. Linear approximation of the dynamical system provides behavior classification through the Jacobian matrix. The eigenvalues are only considered in the regions of interest, consistent with the linear approximation and

the implicated behaviors. In such a way the method can identify five types of behaviors. However, the method can be not useful when significant overlap of motion patterns is present in the scene, or when there is lack of consistent characteristic flow.

Ali and Shah in [16] propose a framework in which Lagrangian Particle Dynamics is used for the segmentation of high density crowd flows and detection of flow instabilities. The authors treat a flow field generated by a moving crowd as an aperiodic dynamical system. Therefore, a grid of particles is overlaid on the flow field in order to monitor the evolution of the particles. Then, a Finite Time Lyapunov Exponent (FTLE) field is used to quantify the amount of particles and to reveal the Lagrangian Coherent Structures (LCS) present in the underlying flow. The LCS divides flow into regions respecting the dynamics of the scene. The changes in the number of flow segments is considered as an instability.

*Classification.* This category includes approaches that exploit classifiers to recognize the behavior of the observed scene. Greenewald and Hero in [17] describe an approach able to learn the normative multi-frame pixel joint distribution and detect deviations from it using a likelihood based approach. The authors use a mean and covariance approach and consider methods of learning the spatio-temporal covariance in the low-sample regime. The approach estimates the covariance using parameter reduction and sparse models. The first method considered is the representation of the covariance as a sum of Kronecker products, which is found to be an accurate approximation in this setting. Then, they consider the sparse a multi-resolution model and apply the Kronecker product methods to it for further parameter reduction, as well as introducing modifications for enhanced efficiency and greater applicability to spatio-temporal covariance matrices.

Idrees et al. in [18] describe an approach to count number of individuals in extremely dense crowds, on a scale not tackled before. Multiple sources of information are used in order to compute an estimation of the number of individuals present in an extremely dense crowd visible in a single image. Due to the common vision problems (e.g. perspective, occlusions, clutters and few pixel per persons), the proposed approach relies on multiple sources such as low confidence head detections, repetition of texture elements (using SIFT), and frequency-domain analysis to estimate counts, along with confidence associated with observing individuals, in an image region. Then, a global consistency constraint on count using Markov Random Field is employed.

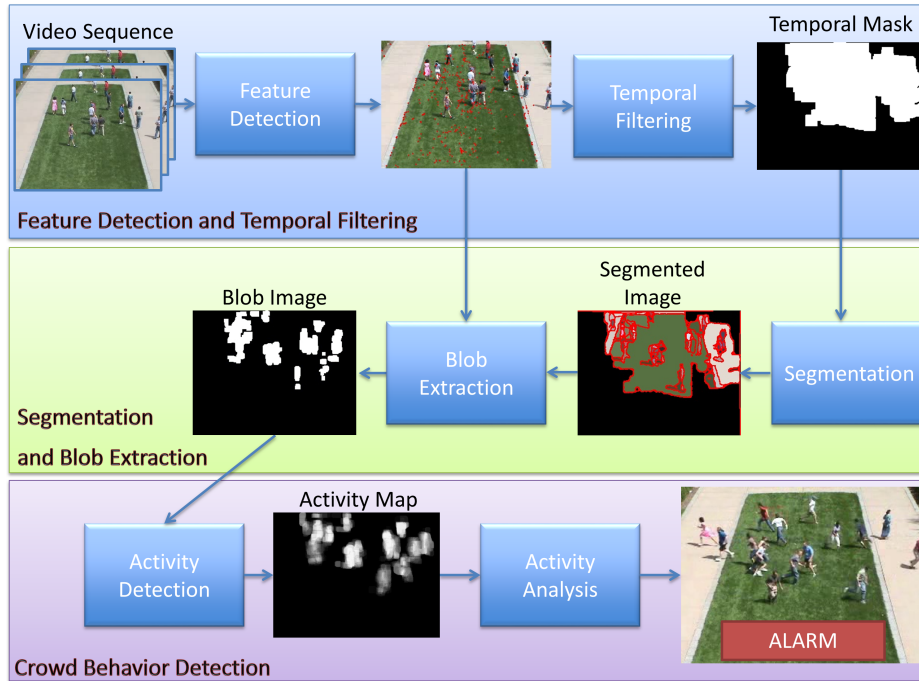


Figure 1: Block diagram of FSCB method.

Moreover, the approach scales well to different densities producing constant error rates across images with diverse count.

In this work, we propose a statistical analysis approach that combines feature detection and image segmentation in order to detect abnormal behaviors in the scene. The proposed method is online and runs in real-time. In particular, two metrics, namely *entropy* and *temporal occupation variation*, are taken into account for detecting abnormal crowd behaviors, without the need of a training phase.

### 3. Feature Tracking and Image Segmentation for Behavior Understanding

In this section, the description of our crowd behavior detection method, called **FSCB**, is provided. FSCB is made of three steps: 1) **F**eature detection and temporal filtering; 2) image **S**egmentation and blob extraction; 3) **C**rowd Behavior detection. The block diagram of FSCB method is shown in Fig. 1

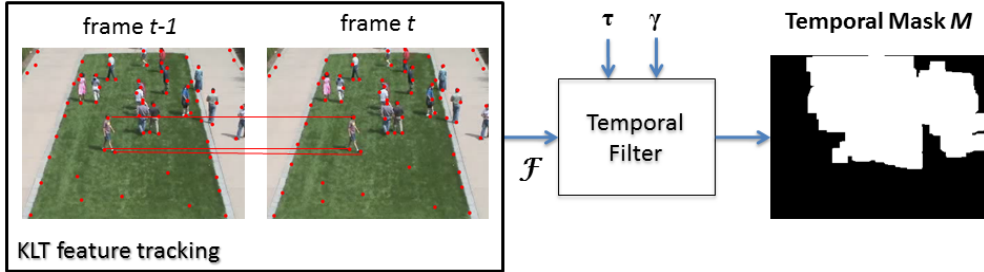


Figure 2: Features in consecutive frames are analyzed to generate a temporal mask  $M$  representing the regions of the scene containing motion.

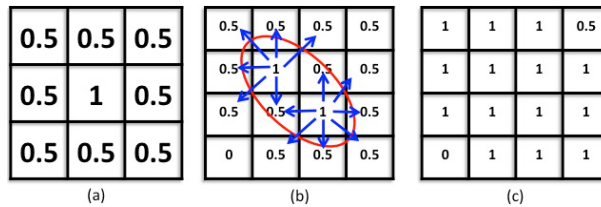


Figure 3: Generation of the probability map: a) The value 1 is assigned for each tracked feature point, while the value 0.5 is assigned to its 8-connected neighbors. b) Once all the features have been considered, if there are adjacent points with value 1, then c) the probabilities of all their neighbors are set to 1 too.

and the details of each step are given in the following.

### 3.1. Feature Detection and Temporal Filtering

The first step of FSCB aims at finding descriptive visual features of the crowd flow in the observed scene. We assume that the following conditions hold in the scene: 1) Brightness constancy, i.e., projection of the same point looks the same in every frame; 2) Small motion, i.e., points do not move very far; 3) Spatial coherence, i.e., points move like their neighbors. The above conditions are usually satisfied in video sequences recorded at 25 frames per second and containing crowded scenes. Given the above assumptions, we decided to exploit the Kanade-Lucas-Tomasi (KLT) feature tracker [19] for detecting and tracking local visual features, instead of using other feature descriptors like Harris corners, SIFT or SURF. Indeed, KLT works very well in situations where distance between images is small, it displays good immunity to tuning parameters, and it has low computational needs [20].

The output of the KLT tracker is a set  $\mathcal{F}$  of couples  $\langle f_i^{t-1}, f_i^t \rangle$ ,  $0 \leq i < n$ ,



of corresponding feature points in two consecutive frames captured at time  $t-1$  and  $t$ , respectively (see Fig. 2). Once  $\mathcal{F}$  has been calculated, a temporal filter is applied in order to create a binary temporal mask  $M$ , containing only the moving points in the scene. To this end, two thresholds are adopted, namely  $\tau$  and  $\gamma$ , to filter out not moving points:  $\tau$  is the length of an history queue, while  $\gamma$  is the minimum velocity value (in pixel per second) to consider a feature point as a moving one.

For each couple  $\langle f_i^{t-1}, f_i^t \rangle \in \mathcal{F}$  a vector  $\mathcal{V} = \{v_1, \dots, v_z\}$ ,  $z \leq \tau$ , is maintained in memory, where  $v_j$ ,  $1 \leq j \leq z$ , represents the velocity, recorded at time  $t-z+j$ , of the feature point  $f_i$ . In particular, the velocity  $v$  of a feature point  $f$  at time  $t$  is calculated as:

$$v = \frac{\sqrt{(f^{t-1}(x) - f^t(x))^2 + (f^{t-1}(y) - f^t(y))^2}}{\text{frame rate in seconds}} \quad (1)$$

At the arrival of every new frame, a set of filtered features  $\mathcal{F}^*$  is obtained by discarding from  $\mathcal{F}$  the features having  $v_z \leq \gamma$ .

Then, a probability grid is used for weighting the motion points  $\mathcal{F}^*$ . The grid has the same size of the input images and it is divided into cells, one for each pixel, and the cells are initialized with the value zero. For each moving point, the 1 value is assigned to the corresponding cell, while the value 0.5 is assigned to all its 8-connected neighbor cells in the grid (see Fig. 3a). After having analyzed all the feature points belonging to  $\mathcal{F}^*$ , the grid is further modified in order to cluster adjacent moving points: If a cell of the grid with value 1 has neighbors with value 1 as well (as in Fig. 3b), then all the cells in their neighborhood are set to 1 also (see Fig. 3c).

Finally, the binary temporal mask  $M$  is generated by considering the cells in the grid with value 1 as white points, and the remaining ones as black points (Fig. 2).  $M$  provides a map of the regions in the image where there are moving pixels. In all our experiments, we set  $\tau$  and  $\gamma$  to 10 frames and 2 pixels per second, respectively.

### 3.2. Image Segmentation and Blob Extraction

The RGB image segmentation is performed by using an approach similar to the one described by Taylor and Cowley in [21]. Firstly, the current RGB frame  $I$  in input is filtered by using the binary mask  $M$ , obtaining a new image  $I^*$ . Then, the image  $I^*$  is segmented according two steps: *Edge Segmentation* and *Delaunay Triangulation*. The former is used for splitting

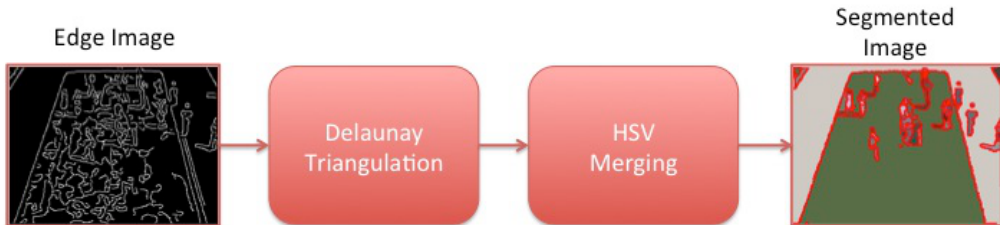


Figure 4: The segmentation step is made of 3 steps: edge detection, Delaunay Triangulation and HSV merging.

the image into local coherent regions, the latter for aggregating homogeneous regions in a global fashion.

$I^*$  is filtered by applying a Gaussian blur filter with a  $3 \times 3$  kernel size. Then, it is converted to grayscale, obtaining an image  $G$ , and the Edge Segmentation procedure begins with a Canny edge extraction, that leads to the creation of an edge image containing the intensity edges in  $G$  (see Fig. 4). The two parameters  $min$  and  $max$  in the Canny algorithm have been set to the values 0.03 and 2.0 respectively, in order to focus on short edges in  $G$ .

The contents of the edge image are then vectorized into connected line segments and used as input for a Delaunay Triangulation procedure, which computes a triangular tessellation of the image.

The Delaunay Triangulation of a point set  $\mathfrak{P}$  is characterized by the empty circumdisk property: no point in  $\mathfrak{P}$  lies in the interior of any triangle's circumscribing disk.

**Definition** [22]. In the context of the finite point set  $\mathfrak{P}$ , a triangle is Delaunay if its vertices are in  $\mathfrak{P}$  and its open circumdisk is empty (i.e., it contains no point in  $\mathfrak{P}$ ). It is worth noting that any number of points in  $\mathfrak{P}$  can lie on a Delaunay triangle's circumcircle. An edge is *Delaunay* if its vertices are in  $\mathfrak{P}$  and it has at least one empty open circumdisk. A Delaunay Triangulation of  $\mathfrak{P}$ , denoted  $Del \mathfrak{P}$ , is a triangulation of  $\mathfrak{P}$  in which every triangle is Delaunay.

Given the connected line segments generated as in [23], the function *Delaunay* from the CGAL<sup>2</sup> library is used to carry out the triangulation. The nodes of the planar triangular graph obtained from the Delaunay Triangulation represent the set of triangles and the edges indicate adjacency relations

---

<sup>2</sup><https://www.cgal.org>

between them, i.e., there is an edge between two nearby triangles.

The triangular graph is segmented using a merging procedure that iteratively finds and merges the two regions with the lowest normalized boundary cost, by considering a predefined association thresholds  $\omega$  (in our experiments, we set  $\omega$  to 0.9). In particular, each one of the triangles in the graph is considered in turn, calculating the average HSV color of all the pixels that lie within its circumcircle. The  $\omega$  threshold is used for measuring the triangle color similarity: If a pair of triangles have a similar normalized HSV value, then they are merged in a single triangle.

An example of the results produced by the image segmentation task is shown in Fig. 4, with the segmentation process carried out on the entire image (instead of focusing only on the moving regions) for better demonstrating the segmentation results. It is worth noting that, in practice, the segmentation process is carried out only on a part of the current frame, denoted by the temporal mask (see Fig. 1).

The extraction of the blobs is performed by applying again the KLT feature tracker, this time on the image  $I^*$ , in order to find the moving blobs. A set  $\mathcal{F}_{blob}$  of couples of corresponding feature points is generated as before. Then, the features are filtered by using Eq. 1, thus obtaining a new set of filtered features  $\mathcal{F}_{blob}^*$ . The set  $\mathcal{F}_{blob}^*$  is re-projected onto the segmented image in order to detect the set  $S$  of moving blobs. A blob is considered as a moving one if its area contains at least a feature point  $f \in \mathcal{F}_{blob}^*$ . In such a way, a binary blob image is obtained (see Fig. 1).

### 3.3. Crowd Behavior Detection

The crowd behavior in the observed scene is detected by carrying out a statistical analysis on the data collected over a temporal window  $w$ . As shown in Fig. 5, given in input a set of binary blob images (Fig. 5a), a 3D-Grid of size  $m \times n \times w$  (Fig. 5b) is used to generate a grayscale activity map (Fig. 5c). The width  $m$  and the height  $n$  of the grid are the same of the input image, while the depth  $w$  corresponds to the length of the temporal window. Then, each voxel  $a$  in the 3D-Grid is set to value 1 if the corresponding pixel  $p$  in the blob image is white. The depth of the voxel  $a$  is represented by a set of values 1, equal to the number of the corresponding white pixels in the blob images, over the time interval  $w$ . In such a way, the temporal persistence of each point  $p$  in the scene is given by the depth of the corresponding voxel  $a$  in the 3D-Grid.

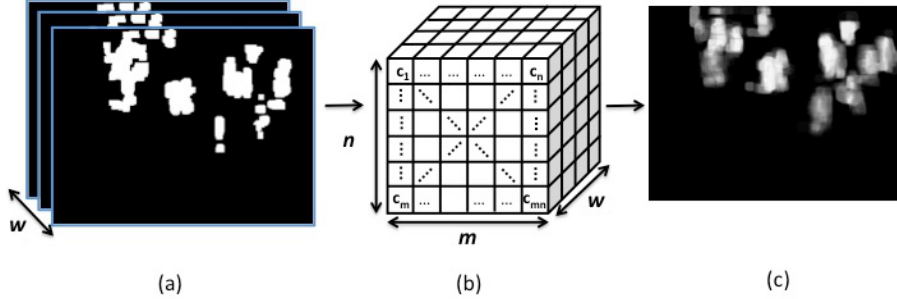


Figure 5: Activity Map Computation: a) A set of blob images is collected over a time window  $w$ . b) A 3D-Grid of size  $m \times n \times w$  is used to record the time persistence of each pixel. c) The activity map is obtained by clustering the data in the 3D-Grid.

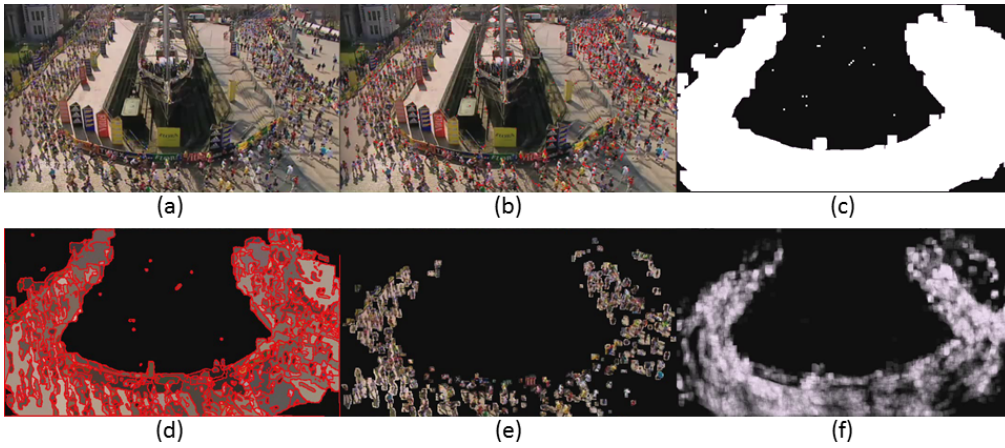


Figure 6: An example of activity map computation. a) The current frame in input (video sequence from [16]). b) Visual feature extraction with KLT. c) Temporal mask. d) Segmented image. e) Moving blob image f) Activity map.

The gray values in the activity map (Fig. 5c) are strictly related to the persistence of the pixel during the time window  $w$ , i.e., a value near 255 in the activity map indicates a point with high activity. In our experiments, the length  $w$  of the temporal window is set to the frame rate value of the video sequence at hand.

Fig. 6 shows all the steps that are performed for obtaining the activity map in a high density crowd scenario. It can be noted that only the part of the image containing a real motion is taken into account.

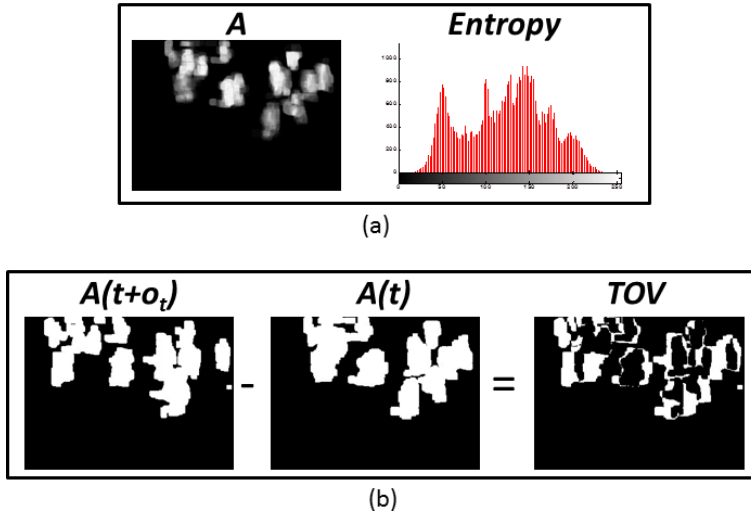


Figure 7: The two measures considered for crowd behavior detection: a) Image entropy. b) Temporal occupancy variation.

Once the activity map is available, it is possible to analyze the trend of two particular measures:

1. *image entropy*;
2. *temporal occupancy variation*.

The image entropy serves for obtaining a measure of the uncertainty in the image values by counting the average amount of information required to encode the image values. The zero order entropy for an image  $I$  is defined as:

$$Entropy(I) = \sum_{i=1}^n p_i \log_2 p_i \quad (2)$$

where  $n$  is the number of separate symbols,  $p_i$  is the frequency of the  $i$ -th pixel in the image, and the result is measured in bits per symbol (pixel value).

Then, by assuming that an infrequent event provides more information than a frequent event [24], it is possible to monitor the instant variation of an image  $I$  in order to detect sudden changes. A threshold  $e_v$  is set as a “sentinel”: If  $Entropy(I(t+1)) - Entropy(I(t)) > e_v$  something of anomalous is happening. An example of image entropy calculation on an activity map  $A$  is shown in Fig. 7a.

The temporal occupancy variation (TOV) takes into account the space occupied by the detected moving blobs over time. Given a temporal threshold  $o_t$ , the TOV is given by:  $TOV = A(t + o_t) - A(t)$ . The value of TOV represents the percentage of image space occupied during a time  $o_t$ . If the value of TOV increases, it means that the scene is changing. We assume that in case of a great variation in the TOV value, an abnormal event is happening. An example of TOV calculation is shown in Fig. 7b.

A discussion about the values for the thresholds  $e_v$  and  $o_t$  is provided in the next section.

## 4. Experimental Evaluation

The experimental results reported in this section are related to the problem of detecting events of interest in crowded scenes. Multiple publicly available video sequences have been selected for quantitatively evaluating the proposed approach and for comparing it with other recent state-of-the-art online approaches.

### 4.1. Data Sets

Four different data sets have been selected for the experiments: UMN [25], PETS 2009 [6], AGORASET [7], and Rome Marathon [26]. Each data set contain one or multiple video sequences and the corresponding ground truth data. Each frame in a video sequence is labeled with a value “normal” or “abnormal”, with “abnormal” meaning that an event of interest is in progress. Ground truth data was already available for the UMN data set, while for the other three data sets we generated the corresponding annotation data, that are available at the following web page <http://www.dis.uniroma1.it/~pennisi/eventdetection.html>. A brief description of the selected data sets is provided in the following.

*UMN Data set.* UMN data set has been collected by the University of Minnesota, USA, and it consists of eleven videos representing escape events. The videos are captured in three different indoor and outdoor scenes, commonly denoted as Lawn, Indoor, and Plaza. Each video starts with a crowd, of about 20 people, that walks in different directions, then an abnormal event causes people to run away. Fig. 8 shows a sample frame for each scene. Ground truth data for UMN are provided by the authors of the data set.

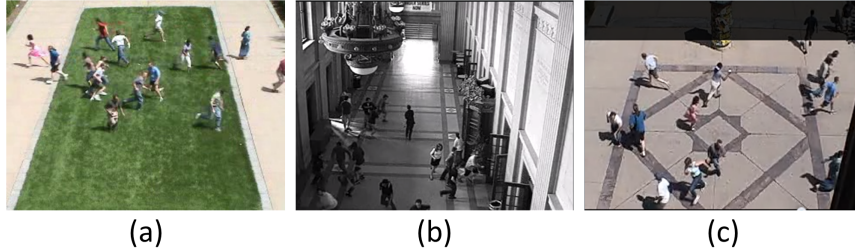


Figure 8: The University of Minnesota (UMN) data set. a) Lawn. b) Indoor. c) Plaza.



Figure 9: The Performance Evaluation of Tracking and Surveillance (PETS) 2009 data set.

*PETS 2009*. The data set has been recorded for the workshop PETS 2009 at Whiteknights Campus, University of Reading, UK. PETS 2009 comprises multi-sensor sequences containing crowd scene scenarios with increasing scene complexity and it is composed by three data sets:

- S1: concerns person count and density estimation;
- S2: addresses people tracking;
- S3: involves flow analysis and event recognition.

In our experiments, we used the S3 data set (Fig. 9) and we manually annotated the sequence creating ground truth data.

*AGORASET*. The AGORASET data set is composed of synthetic scenes representing various crowd simulations. Each video is equipped by different information: ground truth data, the position of the pedestrians, velocity of the flow and a set of MATLAB tools. In AGORASET, seven typical scenes are represented where some crowd behaviors appears. These scenes

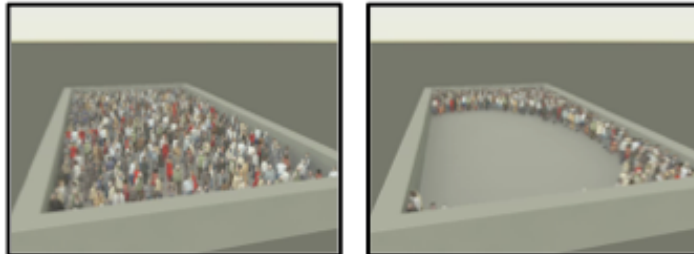


Figure 10: The AGORASET: a data set for crowd video analysis.



Figure 11: The Rome marathon data set.

correspond to an evolution of a human flow in different environments, e.g., an environment with obstacles, an evacuation through a door, etc. In our experiments, we focus on the dispersion scenario (see Fig. 10), where a crowd with about 100 people walks around in a close environment and then moves suddenly to the limit of the environment. In our experiments, we manually annotated the sequence creating ground truth data.

*Rome Marathon.* Since the scarcity of publicly available data set for crowd behavior understanding is an actual problem for the computer vision community, we decided to provide two novel video sequences containing crowded scenes. The data set has been recorded during the 2013 Rome Marathon and it is available for download, together with ground truth data for each video, at [26]. The Rome marathon data set is made of two video sequences representing two different situations: 1) the starting of the marathon and 2) the cleaning of the street. As shown in Fig. 11, the scenes contains thousands of people participating to the marathon.

#### 4.2. Metrics

In order to obtain quantitative results for our FSCB algorithm, we measured the number of frames in the video sequence at hand that are detected



as false positives (FP), true positives (TP), false negatives (FN), and true negatives (TN) with respect to the ground truth data. True positive rate (TPR) and False positive rate (FPR) can be computed with the following formulas:

$$TPR = \frac{TP}{TP+FN}, \quad FPR = \frac{FP}{FP+TN} \quad (3)$$

TPR and FPR can be used for generating a Receiver Operating Characteristics (ROC) curve and for computing the relative Area Under Curvature (AUC).

The area under the ROC is a convenient way of comparing different classification methods. A random classifier has an area of 0.5, while an ideal one has an area of 1. The obtained quantitative results for FSCB on four data sets are provided in the following.

#### 4.3. Quantitative Results

In order to qualitatively evaluate the performance of our FSCB algorithm, we tested the approach generating the ROC curve for each of the above described data sets. All the used ground truth data are publicly available.

It is worth noting that there exist a large variety of offline crowd behavior detection methods that are able to achieve an AUC value near 1 on the considered sequences (e.g., a value of 0.99 is obtained in [27] on UMN). However, such performance are obtained by analyzing the entire video, i.e., having the possibility of exploiting knowledge about events that will happen in the future. This type of analysis is useful in order to obtain a model for different crowd behaviors, but offline analysis can result ineffective for practical use. In the following, we compare our FSCB method only with online state-of-the-art methods.

For the UMN data set a double comparison has been carried out. In the first set of experiment, in order to carry out a fair comparison with published results, the entire data set is considered as a whole video sequence.

The ROC curve generated on the entire UMN sequence (11 videos treated as a single one) is shown in Fig. 12. In particular, the value of  $e_v$  has been varied in the range  $0.1 \leq e_v \leq 0.2$ , while the value of  $o_t$  in the range  $30 \leq o_t \leq 35$ .

Table 1 shows that FSCB achieves better results than the methods relying on pure optical flow (results from [8]) and on a neural network (results from [28]). For FSCB, some false positives are detected due to the anticipated

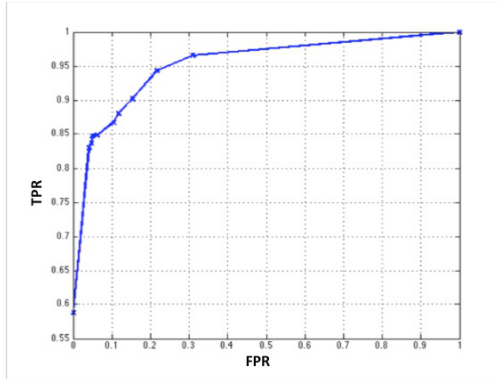


Figure 12: The ROC curve for FSCB on the whole UMN data set. The results are obtained by varying the thresholds  $\tau$  and  $o_t$ .

Table 1: Anomaly detection results on the whole UMN data set. Our approach is compared with other published online methods.

Method	Type	Area under ROC curve (AUC)
Optical Flow [8]	online	0.84
Neural Network [28]	online	0.93
FSCB	online	<b>0.95</b>

detection of the crowd event performed by our approach, with respect to the ground truth data.

The second set of experiments on the UMN data set has been carried out by considering the sequences as divided according to the three different scenarios: Lawn, Indoor and Plaza. Our method has been compared with other two recent online crowd behavior detection methods [29, 30]. Results are reported in Table 2. FSCB performs slightly better than the other two methods on all the three considered sequences.

Along with the well-known UMN data set, three additional video sequences, namely AGORASET, PETS 2009 and Rome Marathon, have been considered for quantitatively evaluating our FSCB method. The results are shown in Table 3. For all the three considered data sets, FSCB is able to achieve good results with an AUC value over 0.90.

FSCB method obtains good detection results on different video sequences, without the need of using a classifier for detecting the crowd behavior in the

Table 2: Anomaly detection results on single sequences of UMN data set. FSCB method is compared with other published methods.

Method	Type	Area under ROC curve (AUC)		
		Lawn	Indoor	Plaza
STCOG [29]	online	0.9362	0.7759	0.9661
COV [30]	online	0.9605	0.8628	0.9746
FSCB	online	<b>0.9641</b>	<b>0.8764</b>	<b>0.9750</b>

Table 3: Anomaly detection results of FSCB method on PETS 2009, AGORASET, and Rome Marathon data sets.

Data set	Area under ROC curve (AUC) for FSCB method
PETS 2009 [6]	0.93
AGORASET [31]	0.94
Rome Marathon [26]	0.96

observed scene. Indeed, FSCB approach is completely online and it does not need any training phase. Qualitative results and the ROC curves for all the considered data sets are shown in Fig. 13.

#### 4.4. Computational Speed

We tested the computational speed of FSCB method in terms of frames per second (FPS). To the best of our knowledge, the computational load for similar approaches in the literature has not been published. The tests have been made by using a commercial notebook with an Intel Core i7 CPU 2.4 GHz 8 GB RAM and a single-threaded C++ implementation of the FSCB algorithm.

From the obtained results it can be noted that, for  $320 \times 240$ , FSCB runs in real-time. When the frame size increases the computational speed for FSCB decreases.

## 5. Conclusions

In this paper, a real-time and online crowd behavior detection algorithm for video sequences is described. The algorithm, called FSCB, is based on a

Table 4: Computational speed for our algorithm on different data sets.

<b>Data set</b>	<b>Image size</b>	<b>Frames per second (FPS)</b>
UMN [25]	$320 \times 240$	20
AGORASET [7]	$640 \times 480$	16
PETS 2009 [6]	$768 \times 576$	11
Rome Marathon [26]	$1920 \times 1080$	5

pipeline made of the following stages: 1) stable features are tracked between frames of the sequence; 2) a temporal mask is extracted; 3) moving blobs are found using segmentation; 4) anomalous events are detected using two measures: instant entropy and temporal occupancy variation.

Quantitative experiments have been conducted on different publicly available data sets: UMN [25], PETS 2009 [6], AGORASET [7]. For PETS 2009 and AGORASET, ground truth data have been produced and made available at the following web page <http://www.dis.uniroma1.it/~pennisi/eventdetection.html>. Furthermore, a novel annotated data set, Rome Marathon [26], containing crowded scenes from the start of a marathon, has been created.

FSCB has been quantitatively compared with other state-of-the-art methods for online crowd event detection. The results of the comparison demonstrate the effectiveness of the proposed approach, that works without the need of a training stage and obtain real-time performance on  $320 \times 240$  frames.

## References

- [1] V. Saligrama, Z. Chen, Video anomaly detection based on local statistical aggregates, in: Computer Vision and Pattern Recognition (CVPR), (2012), pp. 2112–2119.
- [2] C. Brax, L. Niklasson, M. Smedberg, Finding behavioural anomalies in public areas using video surveillance data, in: 11th International Conference on Information Fusion (2008), pp. 1–8.
- [3] J. Junior, S. Mussef, C. Jung, Crowd analysis using computer vision techniques, IEEE Signal Processing Magazine (2010).

- [4] B. Zhan, D. N. Monekosso, P. Remagnino, S. Velastin, L.-Q. Xu, Crowd analysis: a survey, *Machine Vision and Applications* 19 (2008) 345–357.
- [5] M. Rodriguez, J. Sivic, I. Laptev, J.-Y. Audibert, Data-driven crowd analysis in videos, in: *IEEE International Conference on Computer Vision (ICCV)*, 2011, pp. 1235–1242.
- [6] J. Ferryman, A. Shahrokni, Pets 2009 benchmark data, <http://cs.binghamton.edu/~mrldata/pets2009.html>, 2009.
- [7] P. Allain, N. Courty, T. Corpetti, C. Creusot, Agoraset: a dataset for crowd video analysis, <http://www.sites.univ-rennes2.fr/costel/corpetti/agoraset/Site/AGORASET.html>, 2012.
- [8] R. Mehran, A. Oyama, M. Shah, Abnormal crowd behavior detection using social force model, in: *Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 935–942.
- [9] Y. Zhang, L. Qin, H. Yao, Q. Huang, Abnormal crowd behavior detection based on social attribute-aware force model, in: *International Conference on Image Processing (ICIP)*, 2012, pp. 2689–2692.
- [10] M.-C. Chang, N. Krahnstoeber, W. Ge, Probabilistic group-level motion analysis and scenario recognition., in: *International Conference on Computer Vision (ICCV)*, IEEE, 2011, pp. 747–754.
- [11] L. Kratz, K. Nishino, Anomaly detection in extremely crowded scenes using spatio-temporal motion pattern models, 2013 *IEEE Conference on Computer Vision and Pattern Recognition (2009)*.
- [12] H. Fradi, J. Dugelay, Low level crowd analysis using frame-wise normalized feature for people counting, in: *Information Forensics and Security (WIFS)*, 2012 *IEEE International Workshop on*, pp. 246–251.
- [13] S. Srivastava, K. K. Ng, E. J. Delp, Crowd flow estimation using multiple visual features for scenes with changing crowd densities., in: *International Conference on Advanced Video and Signal-Based Surveillance (AVSS)*, IEEE Computer Society, 2011, pp. 60–65.
- [14] L. Wei, W. Xiaojuan, M. Koichi, Z. Hua-An, Foreground detection based on optical flow and background subtract, in: *International Conference on Communications, Circuits and Systems*, pp. 359–362.

- [15] B. Solmaz, B. E. Moore, M. Shah, Identifying behaviors in crowd scenes using stability analysis for dynamical systems, *IEEE Trans. Pattern Anal. Mach. Intell.* (2012).
- [16] S. Ali, M. Shah, A lagrangian particle dynamics approach for crowd flow segmentation and stability analysis, in: *Computer Vision and Pattern Recognition (CVPR)*, 2007, pp. 1–6.
- [17] K. Greenewald, A. Hero, Detection of Anomalous Crowd Behavior Using Spatio-Temporal Multiresolution Model and Kronecker Sum Decompositions, Technical Report, AFRL ATR Center, 2014.
- [18] H. Idrees, I. Saleemi, C. Seibert, M. Shah, Multi-source multi-scale counting in extremely dense crowd images., in: *Computer Vision and Pattern Recognition (CVPR)*, 2013, IEEE, 2013, pp. 2547–2554.
- [19] J. Shi, C. Tomasi, Good features to track, in: *Computer Vision and Pattern Recognition (CVPR)*, 1994, pp. 593–600.
- [20] J. Klippenstein, H. Zhang, Quantitative evaluation of feature extractors for visual slam, in: *Fourth Canadian Conference on Computer and Robot Vision (CRV)*, pp. 157–164.
- [21] C. Taylor, A. Cowley, Parsing indoor scenes using rgb-d imagery, in: *Proceedings of Robotics: Science and Systems*.
- [22] S.-W. Cheng, T. K. Dey, J. Shewchuk, Delaunay mesh generation, CRC Press, 2012.
- [23] Q. Wu, Y. Yu, Two-level image segmentation based on region and edge integration, in: *Digital Image Computing: Techniques and Applications*, pp. 957–966.
- [24] T. Gevers, A. Gijsenij, J. van de Weijer, J. M. Geusebroek, *Color in Computer Vision : Fundamentals and Applications*, Series in Imaging Science and Technology, The Wiley-IS&T, 2012.
- [25] University of Minnesota, Unusual crowd activity data set, <http://mha.cs.umn.edu/Movies/Crowd-Activity-All.avi>, 2006.
- [26] Ro.Co.Co. Lab, Rome marathon data set, <http://www.dis.unrioma1.it/~pennisi/download/romemarathon.zip>, 2014.

- [27] W. Shandong, M. Brian E., S. Mubarak, Chaotic invariants of lagrangian particle trajectories for anomaly detection in crowded scenes, Conference on Computer Vision and Pattern Recognition (2013) 2054–2060.
- [28] Y. Cong, J. Yuan, J. Liu, Sparse reconstruction cost for abnormal event detection., in: CVPR, pp. 3449–3456.
- [29] Y. Shi, Y. Gao, W. Ruili, Real-time abnormal event detection in complicated scenes, in: 20th International Conference on Pattern Recognition (ICPR), pp. 3653–3656.
- [30] T. Wang, J. Chen, H. Snoussi, Online detection of abnormal events in video streams, Journal of Electrical and Computer Engineering (2013) 1–12.
- [31] P. Allain, N. Courty, T. Corpetti, AGORASET: a dataset for crowd video analysis, in: 1st ICPR International Workshop on Pattern Recognition and Crowd Analysis, Tsukuba, Japan, pp. 26–31.

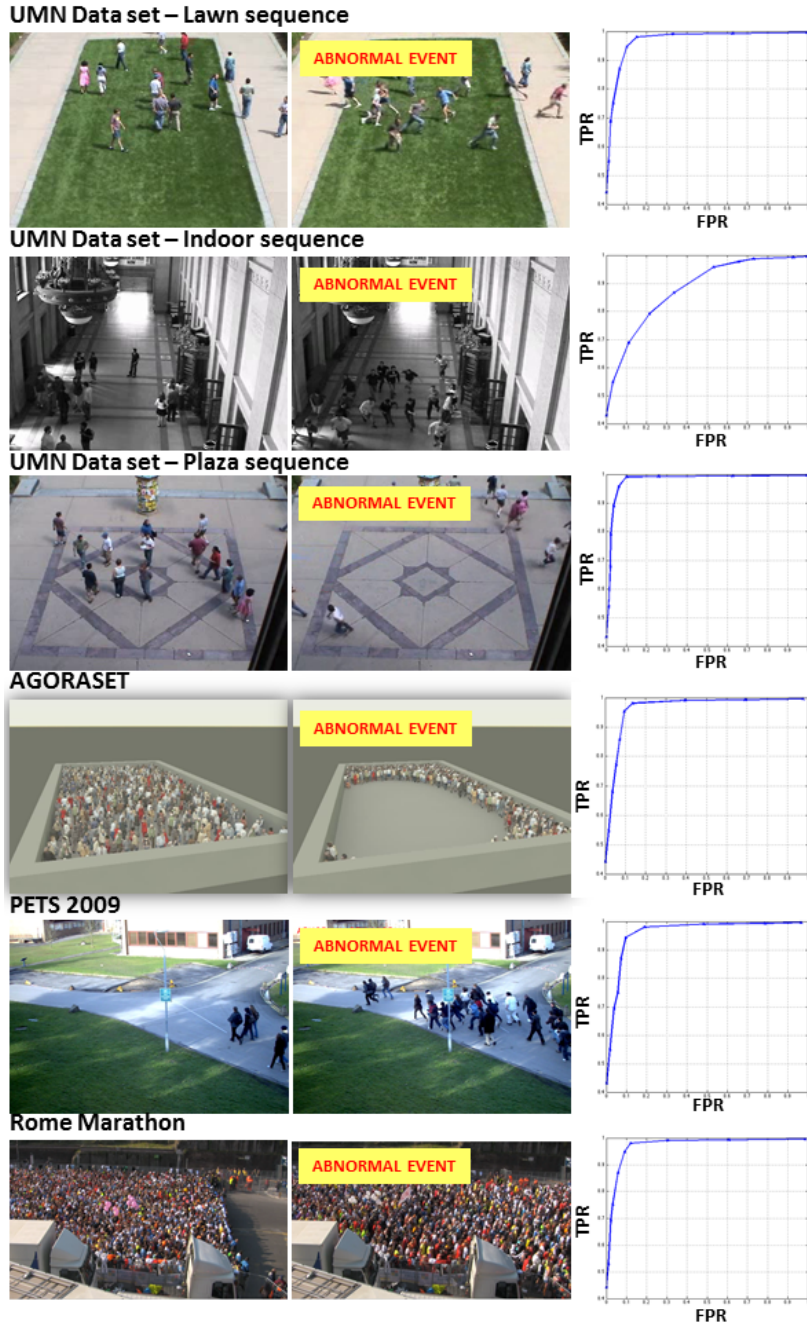


Figure 13: FSCB event detection results. First column: frames without abnormal situations from the considered data sets. Second column: the frames where the abnormal events start. Third column: ROC curves generated by FSCB.