

Detection of Abrupt Changes in Spatial Relationships in Video Sequences

Abdalbassir Abou-Elailah^{1,2(✉)}, Valerie Gouet-Brunet², and Isabelle Bloch¹

¹ LTCI, CNRS, Télécom ParisTech, Université Paris-Saclay, Paris, France
{elailah, isabelle.bloch}@telecom-paristech.fr

² Université Paris-Est, IGN, SRIG, MATIS, 73 avenue de Paris,
94160 Saint Mandé, France
valerie.gouet@ign.fr

Abstract. Detecting unusual events in video sequences is very challenging due to cluttered background, the difficulties of accurate extraction and tracking of moving objects, illumination change, etc. In this work, we focus on detecting strong changes in spatial relationships between moving objects in video sequences, with a limited knowledge of the objects. In this approach, the spatial relationships between two objects of interest are modeled using angle and distance histograms as examples. To evaluate the evolution of the spatial relationships during time, the distances between two angle or distance histograms at two different instants in time are estimated. In addition, a combination approach is proposed to combine the evolution of directional (angle) and metric (distance) relationships. Studying the evolution of the spatial relationships during time allows us to detect the ruptures in such spatial relationships. This study can constitute a promising step toward event detection in video sequences, with few a priori models on the objects.

Keywords: Spatial relationships · Angle histogram · Distances · Fuzzy object representation · Detection of ruptures · Fusion

1 Introduction

In the literature, there are many intelligent video surveillance systems, and each system is dedicated to a specific application, such as sport match analysis, people counting, analysis of personal movements in public shops, behavior recognition in urban environments, drowning detection in swimming pools, etc¹. The VSAM project [35] was probably one of the first projects dedicated to surveillance from video sequences. The goal of ICONS project [18] was to recognize the incidents in video surveillance sequences. The goal of the three projects ADVISOR [2], ETISEO [11] and CareTracker [5] was to analyze record streaming video, in order to recognize events in urban areas and to evaluate scene understanding.

¹ See <http://www.cs.ubc.ca/~lowe/vision.html> for examples of companies and projects on these topics.

The AVITRACK project [3] was applied to the monitoring of airport runways, while the BEWARE project [4] aimed to use dense camera networks for monitoring transport areas (railway stations, metro).

In this context, an increasing attention is paid to “event” detection. In [28], an approach is proposed to detect anomalous events based on learning 2-D trajectories. In [30], a probabilistic model of scene dynamics is proposed for applications such as anomaly detection and improvement of foreground detection. For crowded scenes, tracking moving objects becomes very difficult due to the large number of persons and background clutter. There are many approaches proposed in the literature for abnormal event detection, based on spatio-temporal features. In [19], an unsupervised approach is proposed based on motion contextual anomaly of crowd scenes. In [23], a social force model is used for abnormal crowd behavior detection. In [9], an abnormal event detection framework in crowded scenes is proposed based on spatial and temporal contexts. The same authors proposed in [8] a similar approach based on sparse representations over normal bases. Recently, Hu et al. [16] proposed a local nearest neighbor distance descriptor to detect anomaly regions in video sequences. More recently, the authors in [32] have proposed a video event detection approach based on spatio-temporal path search. It is also applied for walking and running detection.

We adopt a different point of view. We address the question of detecting structural changes or ruptures, which can be seen as a first step for event detection. We propose to use low-level generic primitives and their spatial relationships, and we do not assume a known set of normal situations or behaviors. To our knowledge, the proposed approach is the first one that exploits low-level primitives and spatial relationships in an unsupervised manner to detect ruptures in video. In order to illustrate the interest of spatial relationships, let us consider a person leaving a luggage unattended on the ground. For human beings, it is easy to detect and recognize this kind of event. To learn an intelligent system to detect and recognize this event, one solution is to break down this event into the spatial relationships between the luggage and the person at many points in time. For example, the person holds the luggage at the beginning. If the person leaves the luggage unattended, the spatial relationships between the person and the luggage rapidly changes from very close state to far away state. Thus, detecting ruptures in spatial relationships can be important in detecting and recognizing actions or events in video sequences.

We propose to detect in an unsupervised way strong changes (or ruptures) in spatial relationships in video sequences. This rules out supervised learning-based algorithms which require specific training data. This is useful in all situations where an action or an event can be detected based on such changes or ruptures. Here, we use Harris detector [15], and/or SIFT detector [22] to extract low-level primitives, which are suitable to efficiently detect and track moving objects during time in video sequences [31, 36]. In order to associate features points to objects (to compute the fuzzy representation), the algorithm proposed in [31, 36] can be used. The work presented is considered as a further analysis step after tracking the objects using feature points. Furthermore, we propose a fuzzy representation of the objects, based on their feature points, to improve the

representation of the objects and of the spatial relationships. Then, the structure of the scene is modeled by spatial relationships between different objects using their fuzzy representation. There are several types of spatial relationships: topological relations, metric relations, directional relations, etc. We use directional and metric relationships as an example. More specifically, we consider the angle histogram [24] for its simplicity and reliability, and similarly the distance histogram. In order to study the evolution of the spatial relationships over time and to detect strong changes in the video sequences, we need to measure the changes in the angle or distance histograms during time. Note that this approach differs from methods based on motion detection and analysis, since it considers structural information and the evolving spatial arrangement of the objects in the observed scene. In the literature, many measures have been proposed to measure the distance between two normalized histograms. Here, we propose to adapt these measures to angle histograms, in order to use them in our method. Finally, a criterion is proposed to detect ruptures in the spatial relationships based on distances between angle or distance histograms over time. In addition, a new approach is proposed for combining the distances between angle and distance histograms. The fusion consists in creating a summarized information that represents both the directional and metric spatial relationships. This is a new feature with respect to our preliminary work in [1].

The proposed methods for the fuzzy representation and detection of ruptures in the spatial relationships are described in Sect. 2. Experimental results are shown in Sect. 3 in order to evaluate the performance of the proposed approach. Finally, conclusions and future work are presented in Sect. 4.

2 Rupture Detection Approach

The proposed approach is divided into two main parts. In the first part, our goal is to estimate a fuzzy representation of the objects exploiting only feature points. In the second one, spatial relationships between objects are investigated, using this representation of the objects. Based on the evolution of the spatial relationships during time, strong changes in video sequences are detected.

The fuzzy representation of the objects using the features points is described in Sect. 2.1. Specifically, we study the spatial distribution of the feature points that are extracted using a detector such as Harris or SIFT, for a given object. Feature points can be used to isolate and track objects in video sequences [31, 36]. Thus, we suppose that each moving object is represented by a set of interest points isolated from others with the help of such techniques. Here, we propose two different criteria to represent the objects as regions, exploiting only the feature points. The first one is based on the **depth** of the feature points, by assigning a value to each point based on its centrality with respect to the feature points. The second one assigns a value to each point depending on the **density** of its closest feature points. Finally, the depth and density estimations are combined together, to form a fuzzy representation of the object, where the combined value at each pixel represents the membership degree of this pixel to the object.

This allows reasoning on the feature points or on the fuzzy regions derived from them, without needing a precise segmentation of the objects.

In Sect. 2.2, the computation of the spatial relationships is discussed based on the fuzzy representation of the objects. As an example, we illustrate the concept with the computation of the angle and distance histograms. Then, the existing distances between two normalized histograms are detailed, and the adaptation of these distances to angle histograms is also discussed. Finally, a criterion is defined as the distance between the angle or distance histograms during time, in order to detect ruptures in the spatial relationships.

2.1 Fuzzy Object Representation

In this section, we detail the estimation of the fuzzy representation based on the feature points.

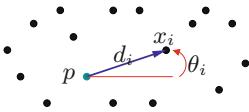


Fig. 1. Feature point distribution for a given object.

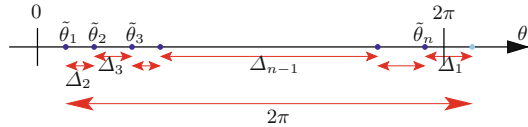


Fig. 2. Sorted angles.

Feature Detection. For a given object, let x_k ($k = 1, 2, \dots, n$) be the detected feature points. For a given pixel p of the object, let px_i denote the line connecting the pixel p and x_i ($i \in [1 \dots n]$), d_i the distance between p and x_i , and θ_i the angle between $\overrightarrow{px_i}$ and the horizontal line as shown in Fig. 1 ($\theta_i \in [0, 2\pi]$).

Distances d_i and angles θ_i are used to estimate depth and density weights for each object based on the x_i . The depth weight is computed using the angles θ_i , and is denoted by dh . The second weight is computed using the distances d_i , and is denoted by dy . Hereafter, their estimations are described, as well as their fusion.

Depth Estimation. In the depth estimation (i.e. centrality), all the feature points are taken into account. Several approaches have been proposed in the literature for depth measures [17], such as simplicial estimation [20], half-space estimation [33], convex-hull peeling estimation [10], L1-depth [34], etc. We propose a new depth measure which is based on the entropy. For each pixel p , the computed angles θ_i are sorted in ascending order as shown in Fig. 2. Let $\tilde{\theta}_i$ ($\tilde{\theta}_j \geq \tilde{\theta}_i$ if $j > i$) be the sorted angles. We define Δ_i as follows:

$$\Delta_i = \begin{cases} (2\pi + \tilde{\theta}_1) - \tilde{\theta}_n & \text{if } i = 1 \\ \tilde{\theta}_i - \tilde{\theta}_{i-1} & \text{if } i \in [2 \dots n] \end{cases} \tag{1}$$

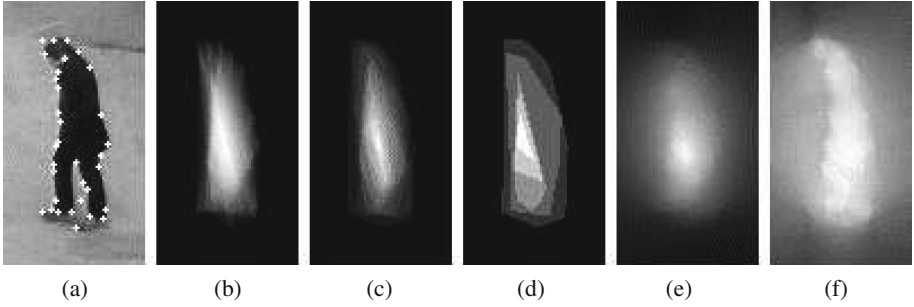


Fig. 3. Depth measures: original object with feature points (a), simplicial estimation [20] (b), half-space estimation [33] (c), convex-hull peeling estimation [10] (d), L1-depth [34] (e), and the proposed depth (f) (image from PETS 2009 database [27]).

Let $p_i = \frac{\Delta_i}{2\pi}$, p_i has two properties: $0 \leq p_i \leq 1$ and $\sum_{i=1}^n p_i = 1$. Thus, p_i can be seen as a discrete probability distribution of the angles. Then, the depth weight is defined as the entropy of this probability distribution:

$$\text{dh}(p) = \frac{1}{n} \sum_{i=1}^n -p_i \log_2 p_i \quad (2)$$

This depth measure can be explained as follows: let us consider a point q inside the object with feature points distributed equitably around it in terms of directions. In this case, we obtain $p_0 = p_1 = \dots = p_n$, and the depth weight of point q is equal to 1 (the highest weight). Otherwise, if the point q is outside the object, the depth weight depends on the angle view (Δ_1 can represent the angle view) and the distribution of the feature points inside the object (p_2, p_3, \dots, p_n). If the angle view becomes smaller and smaller (e.g. the point q is moving away from the object), the depth weight of the point q becomes also smaller accordingly.

Figure 3 shows the representation of several state of the art depth estimations for an object, including our proposal. As we can see, the entropy depth can better represent the shape of the object than the existing depth measures. In terms of computation time, the L1-depth and the proposed depth are the most efficient ones compared to other measures. Our experimental tests showed that the choice of a particular depth measure has a limited impact on the detection of the rupture. However, the entropy depth measure may present a significant enhancement compared to other depth measures, in the applications that need a precise shape estimation, to describe fine relationships, for example when objects meet.

Density Estimation. For density estimation, for a given pixel inside the object, only the neighbor feature points are taken into consideration (feature points within a certain distance r , or k closest feature points). Thus, the distances d_i that are lower than a certain distance r are taken into account to compute the density weight for the pixel p as follows:

$$dy(p) = \sum_{i=1}^M \left(1 - \frac{d_i}{r}\right), \text{ where } d_i \leq r \quad (3)$$

where M is the number of points inside the circle of radius r . This radius can be estimated automatically and online, based on statistics on the distances between points, in order to be adapted to the scale of the object. Figure 4(c) shows a representation of the density estimation.

Fusion of Depth and Density Estimations. We present a combination approach to fuse the two estimations obtained from depth and density of the feature points. For the sake of optimization, the pixels q that are taken into consideration for the fusion are defined as follows: $dy(q) > 0$ or $dh(q) > th$, where th is a given threshold. The obtained estimation of the object is referred to as “fuzzy representation”.

Here, the z-score [6] is applied on the two estimations, in order to make them comparable. The z-score is the most commonly used normalization process. It converts all estimations to a common scale with an average of **zero** and a standard deviation of **one**. It is defined as follows: $Z = (X - \bar{M})/(\sigma)$, where \bar{M} and σ represent the average and the standard deviation of the X estimation, respectively. Let Z^{dh} and Z^{dy} be the depth and density estimations respectively, after applying the z-score normalization.

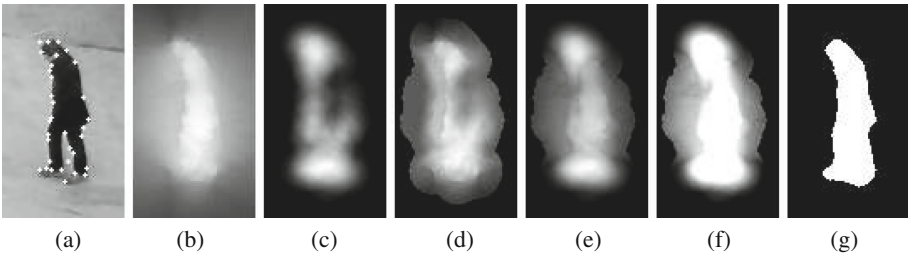


Fig. 4. Original object with the feature points (a), depth estimation (b), density estimation (c), fusion using min operator (d), fusion using max operator (e), fusion using Eq. 4 (f), and the object segmented precisely GT (g).

The obtained fuzzy representation, using different fusion operators, are compared with a Ground Truth (GT) where the objects are segmented precisely (see Sect. 3 for details, and an example in Fig. 4(g)). The combination approach which gives the best performance consists in using the two operators min and max together as defined in the following expression:

$$F(p) = \min \left(\max \left(Z^{dh}(p), Z^{dy}(p) \right), \hat{\sigma} \right) \quad (4)$$

where $\hat{\sigma} = \frac{1}{2th}$. Then, F is normalized using Min-Max scaling [14] to obtain the membership function μ_F which varies in $[0, 1]$. This fusion can be explained as

follows: when Z^{dh} (or Z^{dy}) is greater than $\hat{\sigma}$, the membership value $\mu_F(p)$ is equal to 1. Otherwise, $\mu_F(p)$ is less than 1 according to the maximum between them. As an example, Fig. 4 shows different fuzzy representations of the object using min operator, max operator, and Eq. 4 for the fusion. As we can see, the last fusion approach shows the best fuzzy representation of the object according to the ground truth. The obtained fuzzy representations are used to compute the spatial relationships.

2.2 Spatial Relationships and Rupture Detection

Here, the goal is to estimate the spatial relationships between two objects based on their fuzzy representation. The angle [24] and distance histograms are selected as examples to model the spatial relationships. It is important to note that the proposed method also applies to other types of spatial relationships.

Angle Histogram. Given two fuzzy regions $A = \{(a_i, \mu_A(a_i)), i = 1, \dots, n\}$ and $B = \{(b_j, \mu_B(b_j)), j = 1, \dots, m\}$, where a_i and b_j are the elements of A and B , and μ_A and μ_B represent their membership functions respectively, for all possible pairs $\{(a_i, b_j), a_i \in A \text{ and } b_j \in B\}$, the angle θ_{ij} between a_i and b_j is computed, and a coefficient $\mu_\Theta(\theta_{ij}) = \mu_A(a_i) \times \mu_B(b_j)$ is derived. For a given direction α , all the coefficients of the angles that are equal to α are accumulated as follows:

$$h^\alpha = \sum_{\theta_{ij}=\alpha, i=1, \dots, n, j=1, \dots, m} \mu_\Theta(\theta_{ij}) \quad (5)$$

Finally, $h = \{(\alpha, h^\alpha), \alpha \in [0, 2\pi]\}$ is the angle histogram. In our case, the histogram can be seen as an estimate of the probability distribution of the angles. Thus, the obtained histogram is normalized to display frequencies of the existed angles with the total area equaling 1. It is normalized by dividing each value by the sum $R_h = \sum_{\alpha \in [0, 2\pi]} h^\alpha$, instead of normalizing by the maximum value (which would correspond to a possibilistic interpretation).

When the objects are represented sparsely by feature points, then $\mu_A(a_i) = 1$ and $\mu_B(b_j) = 1$ (where a_i and b_j represent the feature points on the objects A and B respectively), and the same approach is used to compute the angle histogram between the two sparse objects A and B .

Distance Histogram. In this case, all the distances d_{ij} between a_i ($i = 1, \dots, n$) and b_j ($j = 1, \dots, m$) are computed. Based on these distances, the distance histogram is formulated in the same way as the angle histogram:

$$h^l = \sum_{d_{ij}=l, i=1, \dots, n, j=1, \dots, m} \mu_L(d_{ij}) \quad (6)$$

where $\mu_L(d_{ij}) = \mu_A(a_i) \times \mu_B(b_j)$ and l represents a given distance value. The obtained histogram is normalized such that the sum of all bins is equal to 1.

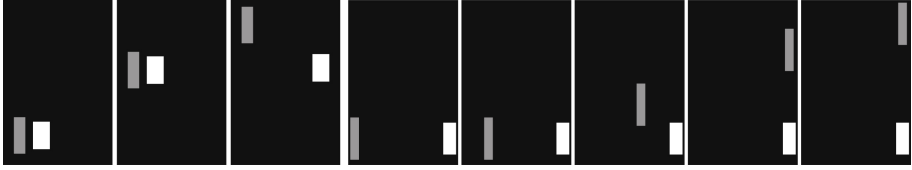
Comparison of Spatial Relationships. There are two main approaches to estimate distances between histograms. The first approach is known as bin-to-bin

distances such as L_1 and L_2 norms. The second one is called cross-bin distances; it is more robust and discriminative since it takes the distance on the support of the distributions into account. Note that the bin-to-bin distances may be seen as particular cases of the cross-bin distances. Several distances based on cross-bin distances, such as Quadratic-Form (QF) distance [13], Earth Mover's Distance (EMD) [29], Quadratic-Chi (QC) histogram distance [25], have been proposed in the literature. We have tested these three distances on different examples, and experiments showed that they were well adapted to angle histograms. Finally, the QF distance was used in our experiments to assess the distance between the angle or distance histograms during time, because of its simplicity. It is defined as follows: $d(h_1, h_2) = \sqrt{ZSZ^T}$, where $Z = h_1 - h_2$ and $S = \{s_{ij}\}$ is the bin-similarity matrix. This distance is commonly used for normalized histograms (the distance histogram for example). Here, we propose an approach to adapt it to the case of angle histograms just by adjusting the elements of the similarity matrix S . We consider that the two histograms h_1 and h_2 defined on $[0, 2\pi]$ consist of k bins B_i . Usually, for a distribution on the real line, the distance between B_i and B_j is defined as follows: $x_{ij} = |B_i - B_j|$, where $1 \leq i \leq k$ and $1 \leq j \leq k$. However, in the case of angle histograms, the distance between B_i and B_j is defined as follows: $x_{ij}^c = \min(x_{ij}, 2\pi - x_{ij})$ to account for the periodicity on $[0, 2\pi]$. Thus, the elements of the matrix S are simply defined, in the case of angle histograms, using x_{ij}^c instead of x_{ij} as follows:

$$s_{ij} = 1 - \frac{x_{ij}^c}{\max_{i,j}(x_{ij}^c)} \quad (7)$$

Criterion for Rupture Detection. Based on the fuzzy representation of the objects exploiting only the feature points, the angle or distance histogram h between two different objects is computed. Let f_i ($i = 0, 1, \dots, N - 1$) be the frames of the video sequences, and h_i be the computed angle or distance histogram between the objects A and B in frame f_i . We define $y(i) = d(h_i, h_{i+1})$ for each $i = 0, 1, \dots, N - 1$. This function describes the evolution of the angle or distance histograms over time. If a strong change in the spatial relationships occurs at instant R ($R < N$), where R denotes the instant of rupture, this means that the angle or distance histogram h_R effectively changes compared to previous angle or distance histograms ($h_i, i < R$). A rupture is detected according to the following criterion W: $\forall i < R - 1, y(R - 1) - y(i) > t$, and t is a threshold value. Thus, the instant of rupture R can be effectively detected from the analysis of the function y .

Here, in order to clearly show the instant of ruptures in the spatial relationships and remove noise, we also show the evolution of the function y filtered by a Gaussian derivative, denoted by g , instead of a simple finite difference. This filter can remove noise and the function g effectively exhibits the instant of the strong changes in the spatial relationships using a threshold approach. This approach is particularly well suited for abrupt changes, leading to clear peaks in the function g , that are then easy to detect (a simple threshold can be sufficient). For slower changes, a multiscale approach can be useful to detect more spread peaks.



(a) Frames number 1, 30, and 50 of SE 1. (b) Frames number 45, 55, 74, 95, and 105 of SE 2.



(c) Frames number 450, 462, and 468 of RE 1 selected from PETS 2009. (d) Frames number 595, 630, 670, and 700 of RE 2 selected from PETS 2009.

Fig. 5. Events SE 1 (a), SE 2 (b), RE 1 (c) and RE 2 (d).

Fusion of Directional (angle) and Metric (distance) Evolutions. To distinguish between two functions that are derived from angle and distance histogram, let y^θ and y^d be the functions that represent the evolution of directional (angle) and metric (distance) spatial relationships during time respectively. The goal of this study is to combine the two functions y^θ and y^d in an efficient way, in order to produce a unique function y^u , which allows us to detect the strong changes in both directional and metric relationships, at the same time. Thus, if a rupture occurs in at least one of them (y^θ and y^d), this rupture must be efficiently detected using the function y^u .

To combine the two evolutions, at each instant time k , the two distances $y^\theta(k)$ and $y^d(k)$ are extracted and used to provide a single point $p_k(x_k, y_k)$ in \mathbf{R}^2 , defined as follows :

$$\begin{cases} x_k = (d_0 + y^d(k)) \cos(y^\theta(k)) \\ y_k = (d_0 + y^d(k)) \sin(y^\theta(k)) \end{cases} \quad (8)$$

where d_0 is a constant, to account for the variation of the distance $y^d(k)$ when the value of $y^d(k)$ is very small (e.g. close to 0). Furthermore, we define a single function y^u using the points p_k ($k = 1, \dots, n$), by computing the distances between two consecutive points p_k and p_{k+1} , over time. The value of the function y^u at instant time k is computed as follows:

$$y^u(k) = \sqrt{(x_{k+1} - x_k)^2 + (y_{k+1} - y_k)^2} \quad (9)$$

this function y^u can be used to detect the ruptures in both directional and metric relationships, using the approach described above.



(a) Frames number 1, 5, 10, and (b) Frames number 1955, 2010, 2060, and 2100 of RE 3 selected from SE 3.

Fig. 6. Events SE 3 (a) and RE 3 (b).

3 Experiments and Evaluations

To evaluate the performance of the proposed approach, we created some synthetic events (illustrated in Fig. 5(a) and (b)), and also used a variety of events selected from the PETS 2009 datasets [27] (illustrated in Fig. 5(c) and (d)). Here, we call “event”, some frames that contain a rupture in the spatial behavior. The results of the proposed fuzzy representation are also compared to classical segmentation approaches: a binary segmentation approach [7] and an approach using differences between the background and the actual frame. Then, morphological operations are carried out to remove small objects and fill holes. The last one is used as ground truth (GT) because it produces very precise segmentations.

A synthetic event and an event selected from PETS 2006 dataset [26], displayed in Fig. 6, are used to illustrate the proposed approach using the distance histogram. To associate feature points to objects, here we simply consider the points included in the bounding boxes associated with objects available in the PETS 2009 dataset.

3.1 Parameters Tuning

In this section, some results are detailed concerning the tuning of the parameters that are used in the proposed approach. Specifically, we discuss the estimation of the radius r , which is used in the computation of the density estimation. Then, some results are shown for different values of the threshold th , which is used in the combination of depth and density estimations. Finally, we show the effect of the number of bins on the computation of the distance between two angle histograms.

r Parameter. Fig. 7 shows different estimations of the radius r (normalized) during time. First, all the possible distances d_{ij} among the feature points are computed. The mean, median, and maximum of these distances are computed, as shown in the figure (three first curves). Then, Delaunay triangulation is applied on the feature points, and two other estimations of the radius r are computed, as the mean and median of the lengths of the triangle edges (fourth and fifth curves). Finally, as in [21], the median of all radius of the circumscribed circle around the Delaunay triangles provides the last estimation (last curve). As we can see,

the maximum of the distances (third curve) gives the most robust and stable estimation during time. Other experiments on different objects show the same result. Thus, the expression

$$r = \max_{i=1, \dots, n, j=i, \dots, m} \frac{d_{ij}}{6} \quad (10)$$

is adopted to estimate the radius r for the density estimation.

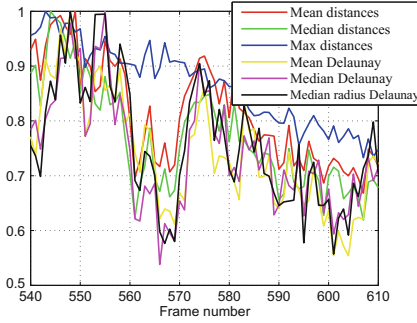


Fig. 7. Different estimations of the radius r based on the feature points.

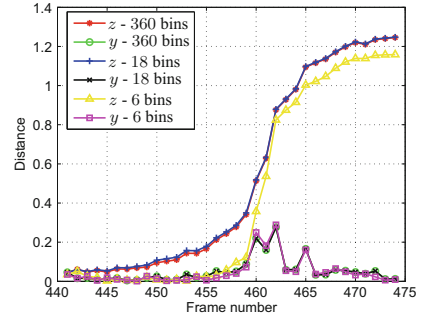


Fig. 8. The functions y and z over time using various number of bins.

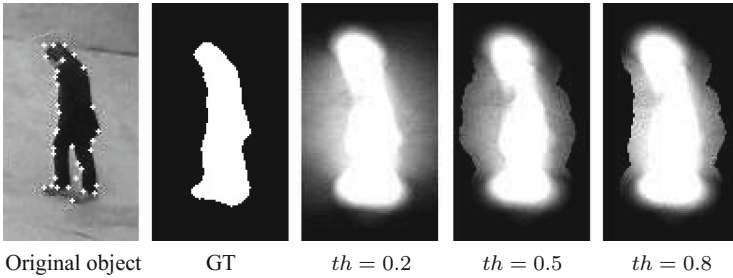


Fig. 9. Original object with the feature points, GT of the object, and fuzzy representations of the object for th equal to 0.2, 0.5, and 0.8 respectively.

th Parameter. In the fusion of depth and density estimations, a threshold th is used. Figure 9 shows the original object with the feature points, the ground truth (GT) of the object, and the fuzzy representation (FR) of the object for different values of th . As we can see, the proposed fusion approach is quite robust to the variation of the used threshold th . In the paper, a value of th equal to 0.5 is used in the combination of depth and density estimations.

Number of Bins. In this section, we study the effect of the number of bins (quantification) on the distance between two angle histograms. We defined the

function y as the distance between two successive angle histograms in frames f_i and f_{i+1} . Here, we also define $z(i) = d(h_0, h_i)$ for $i = 0, 1, \dots, N - 1$, i.e. the distance to the histogram in the initial frame, to consider strong changes in the angle histograms. Figure 8 shows the evolution of the two functions y and z , for numbers of bins of 360, 18, and 6. As we can see, there is almost no difference between 360 and 18 bins, for the two functions. For a number of bins equal to 6, there is a difference compared to 360 and 16 bins for the function z . For the function y , the three curves are almost the same. Thus, the used distance between two angle histograms is robust to the variation of the number of bins.

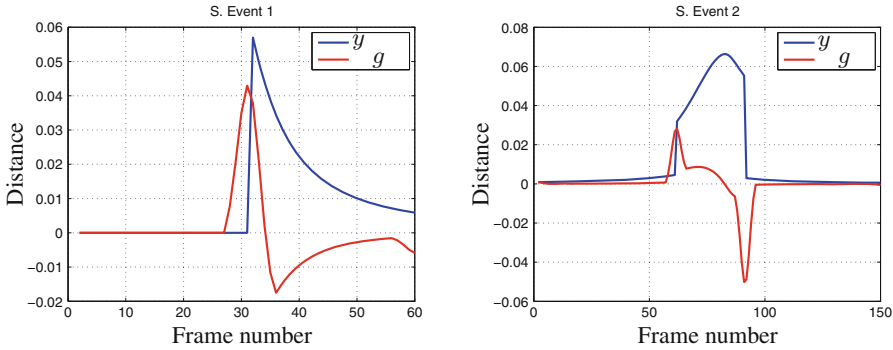


Fig. 10. Functions y and g for events SE 1 (left) and 2 (right), computed from angle histograms.

3.2 Ruptures in Spatial Relationships

We now illustrate how the analysis of the distances between histograms allows us to detect ruptures in spatial relations, both for orientation and distances.

Angle Histogram. Three snapshots of the first synthetic event (SE) are shown in Fig. 5(a) (two objects moving together and then separately). In this case, there is a rupture in the directional spatial relationships, when the two objects diverge. Figure 5(b) shows five snapshots of the second SE. In this event, the object B moves towards the object A (fixed) from the left to the right. Then, the object B changes of direction (frame 74), and when the object B becomes above the object A , it goes towards the top.

Figure 10 shows the functions y and g during time for the two events SE 1 and 2. For the event SE 1, the function y shows a strong variation at frame number 31. At this instant, there is the rupture in the spatial relationships (the two objects begin to separate). Using the evolution of g over time, the instant of the rupture can be detected by applying a threshold (a threshold of 0.02 can be used to detect the instants of rupture for the SE). For the second SE, we can see two strong variations in the function y ; the first strong variation (frame 60) occurs when B

changes of direction with respect to A , the second strong variation (frame 90) occurs when B becomes above A and changes its direction towards the top. The function g clearly shows the two strong variations. Thus, the proposed method can efficiently detect the instants of ruptures in the spatial relationships. Other SE were created and tested using the proposed approach, and similar results were obtained.

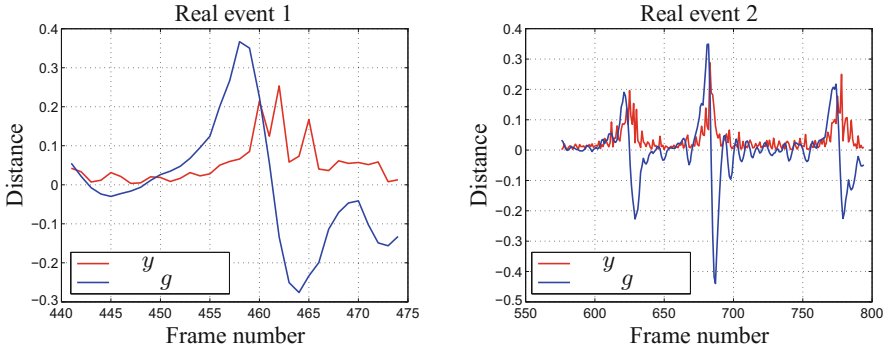


Fig. 11. Functions y and g over time, using the proposed fuzzy representation, for the events RE 1 (left) and RE 2 (right), computed from angle histograms.

Let us now evaluate the proposed detection of ruptures in the spatial relationships in the presence of noise (deformation of objects, etc.) in real events. For the real event (RE) 1 (Fig. 5(c)), the two persons converge then diverge. Figure 11 (left) shows the functions y and g over time using the proposed fuzzy representation, for the event RE 1. Two ruptures in the directional spatial relationships exist in this event. The first one is when the two persons meet, and the second rupture when the two persons separate. It is clear that the two instants of the ruptures can be efficiently detected using the evolution of g (a threshold of 0.2 can be used to detect the instants of ruptures for the RE). In the event RE 2 (Fig. 5(d)), the two persons (surrounded by white and blue bounding boxes) converge and diverge several times. In Fig. 11 (right), we show the functions y and g over time, using the fuzzy representation of the objects, for the event RE 2. All the ruptures in the directional spatial relationships can be efficiently detected using the function g .

Distance Histogram. Four snapshots of the third synthetic event are shown in Fig. 6(a). At the beginning of this event, the two objects diverge at a speed of 5 pixels/frame, and at a given instant (precisely at frame 10), the speed of the two objects becomes 10 pixels/frame. Thus, the velocity of the objects is suddenly increased. Figure 6(b) shows four snapshots of the third real event selected from PETS 2006. In this event, the luggage is attended to by the owner for a moment, and then the person leaves the place and goes away.

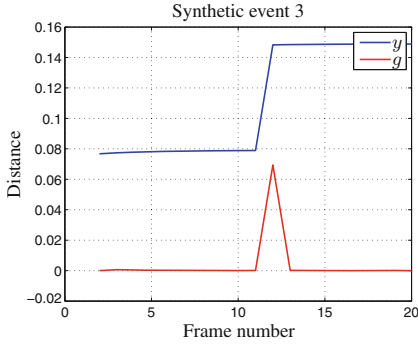


Fig. 12. Functions y and g over time, using the proposed fuzzy representation, for the event SE 3, computed from distance histograms.

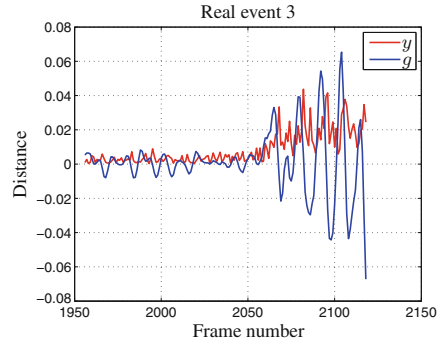


Fig. 13. Functions y and g over time, using the proposed fuzzy representation, for the event RE 3, computed from distance histograms.

In Fig. 12, the functions y and g during time for the event SE 3 are shown. As we can see, the function y shows a strong variation at frame number 10, when the velocity of the objects changes. At this instant, a rupture in the metric spatial relationships is detected, using the evolution of g over time.

In the presence of noise, we show in Fig. 13 the functions y and g during time for the third real event. When the person leaves the place and goes away, we can see a strong change in the function y . By analyzing the obtained results, the instant of rupture in the metric spatial relationships can be detected. These results can be used to indicate events occurring in the video sequences, such as escaping in Fig. 6(a) and Left-Luggage in Fig. 6(b).

Fusion of Angle and Distance Histograms. To evaluate the performance of the proposed approach for the fusion of directional and metric information, we create a synthetic event that contains many ruptures, in directional and metric spatial relationships, during time. In Fig. 14, we show the obtained functions y^θ and y^d over time. Note that y^θ and y^d represent the evolution of spatial relationships using the angle histogram and the distance histogram respectively. As we can see, there are many directional and metric ruptures in these functions. On the right side, we show the results of the combination of the two functions y^θ and y^d using a naive method (in this approach, the fusion consists in averaging the two functions) and the proposed approach. As we can observe, the proposed approach shows clearly the instant of ruptures in both directional and metric spatial relationships, and provides higher values than the naive fusion for most of the ruptures.

As we can see, when a rupture occurs in both functions y^θ and y^d , it is clearly shown in the function y^u (see instants 10 and 94). In addition, the last rupture in the function y^d (at instant 100) can be efficiently detected using the function

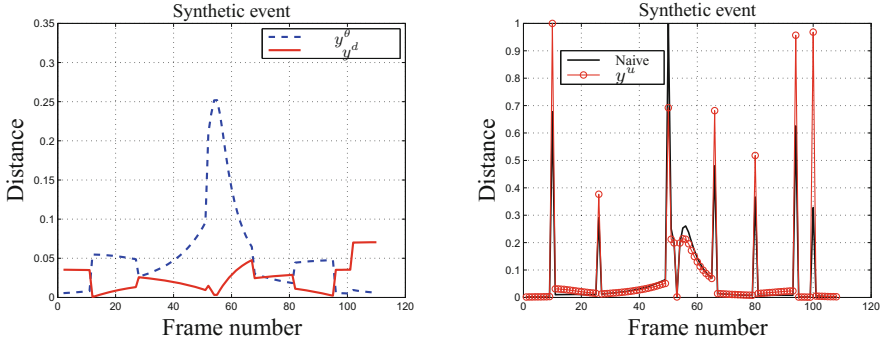


Fig. 14. Functions y^θ and y^δ over time, for a synthetic event (left) and the combination of the functions y^θ and y^δ using the naive fusion and the proposed combination (right).

y^u , even if there is no strong change in the function y^θ at this instant. Similarly, the rupture at instant 50 in the function y^θ can be efficiently detected using y^u , even if at this instant, no rupture is shown in the function y^δ . Thus, the function y^u can show clearly the ruptures that occur in at least one of the functions y^θ and y^δ .

3.3 Impact of Object Representation

Here, we show the importance of the fuzzy representation based on a simple feature points representation. Two feature detectors, Harris and SIFT, are tested. Figure 15 illustrates the function y during time (computed here from angle histograms) for different representations of the objects, for RE 1. The Harris and SIFT features are directly used to estimate the spatial relationships between the two objects and to compute the function y (red and green curves in the figure). In addition, we show in the same figure the evolution of the function y computed on the fuzzy representation of the objects using the Harris and SIFT features (blue and black curves in the figure). As we can see, the evolution of the function y obtained from the fuzzy representation of the objects using the SIFT features (black curve) can significantly reduce the variation of the distance (i.e. less amplitude of the curve) on areas when there is no rupture in the spatial relationships (see Fig. 15, frames 440 to 456) with respect to the SIFT features without computing the fuzzy representation. Thus, the proposed fuzzy representation of the objects before computing the spatial relationships can improve the robustness of the detection of ruptures, based on the observation that SIFT features are more noisy across frames than Harris features in this sequence.

However, noise is present in the function y for all object representations. Assuming that the function y has additive Gaussian noise, the algorithm proposed by Garcia [12] is used to estimate the variance of the noise of the function y , for the different object representations: Harris features, fuzzy representation

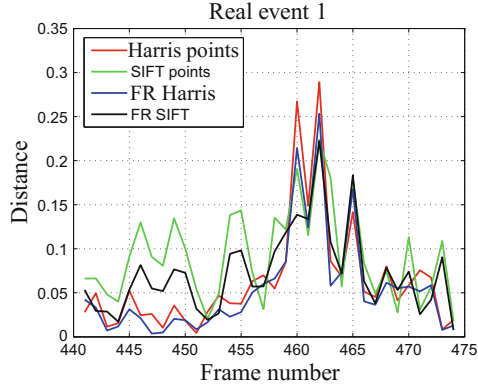


Fig. 15. Function y over time, computed from angle histograms, for different estimations of the objects: Harris features, SIFT features, fuzzy representation (FR) of the objects using Harris features (FR Harris) and SIFT features (FR SIFT), for RE 1.

of the objects using Harris features (FR Harris), SIFT features, fuzzy representation of the objects using SIFT features (FR SIFT), the binary segmentation using Mean-Shift algorithm [7] and GT.

Table 1. Estimated variance of the noise ($\times 10^{-4}$) [12] in the function y , for different object representations, for RE 1 and 2.

Event	Harris	FR Harris	SIFT	FR SIFT	Mean-Shift	GT
RE 1	13	12	27	10	31	12
RE 2	7.7	5.48	8.9	7	31	5.4

Table 1 shows the variance of the noise in the function y , for the different object representations, for the two events RE 1 and 2. It is clear that the proposed fuzzy representation significantly reduces the variance of the noise, which becomes close to the one of the GT. Especially, for SIFT features, the variance of the noise reduces from 27 to 10 for RE 1, and from 8.9 to 7 for RE 2. In addition, the variance of the noise of the proposed object representation is significantly less than the one of the binary segmentation using Mean-Shift algorithm.

4 Conclusion

In this paper, a new method was proposed to detect strong changes in spatial relationships in video sequences. Specifically, new approaches have been proposed to compute depth and density estimations, based on feature points, as well as fuzzy representations of the objects by combining depth and density

estimations. Exploiting the fuzzy representations of the objects, the angle and distance histograms are computed. Then, the distance between the angle or distance histograms is estimated during time. Based on these distances, a criterion is defined in order to detect the significant changes in the spatial relationships. A new approach has been also proposed to combine directional and metric spatial relationships. The proposed method shows good performances in detecting ruptures in the spatial relationships for both synthetic and real video sequences.

Future work will focus on further improvement of the proposed method in order to detect other kinds of ruptures, and investigating the use of spatio-temporal relationships. Besides, we will investigate multi-time scale analysis, in order to better detect events that take more time to happen. In addition, proposing a complete event detection framework based on spatial relationships as discriminative features seems to be promising.

Acknowledgements. This research is part of French ANR project DESCRIBE “Online event detection in video sequences using structural and Bayesian approaches”.

References

1. Abou-Elailah, A., Gouet-Brunet, V., Bloch, I.: Detection of ruptures in spatial relationships in video sequences. In: International Conference on Pattern Recognition Applications and Methods (ICPRAM), pp. 110–120 (2015)
2. Advisor. Advisor Project (2000). <http://www-sop.inria.fr/orion/ADVISOR/>
3. Avitrackr. Avitrackr Project (2004). <http://www-sop.inria.fr/members/Francois.Bremond/topicsText/avitrackrProject.html>
4. Beware. Beware Project (2007). <http://www.eecs.qmul.ac.uk/~sgg/BEWARE/>
5. Caretaker. Caretaker Project (2006). <http://www-sop.inria.fr/members/Francois.Bremond/topicsText/caretakerProject.htm>
6. Carroll, S.R., Carroll, D.J.: Statistics made simple for school leaders: data-driven decision making. R&L Education (2002)
7. Comanicu, D., Meer, P.: Mean shift: a robust approach toward feature space analysis. *IEEE Trans. Pattern Anal. Mach. Intell.* **24**(5), 603–619 (2002)
8. Cong, Y., Yuan, J., Liu, J.: Abnormal event detection in crowded scenes using sparse representation. *Pattern Recogn.* **46**(7), 1851–1864 (2013)
9. Cong, Y., Yuan, J., Tang, Y.: Video anomaly search in crowded scenes via spatio-temporal motion context. *IEEE Trans. Inf. Forensics Secur.* **8**(10), 1590–1599 (2013)
10. Eddy, W.: Convex hull peeling. In: COMPSTAT, pp. 42–47 (1982)
11. Etiseo (2004). <http://www-sop.inria.fr/orion/ETISEO/>
12. Garcia, D.: Robust smoothing of gridded data in one and higher dimensions with missing values. *Comput. Stat. Data Anal.* **54**(4), 1167–1178 (2010)
13. Hafner, J., Sawhney, H., Equitz, W., Flickner, M., Niblack, W.: Efficient color histogram indexing for quadratic form distance functions. *IEEE Trans. Pattern Anal. Mach. Intell.* **17**(7), 729–736 (1995)
14. Han, J., Kamber, M., Pei, J.: *Data Mining: Concepts and Techniques*. Morgan Kaufmann, San Francisco (2006)
15. Harris, C.G., Stephens, M.J.: A combined corner and edge detector. In: Fourth Alvey Vision Conference, pp. 147–151 (1988)

16. Hu, X., Hu, S., Zhang, X., Zhang, H., Luo, L.: Anomaly detection based on local nearest neighbor distance descriptor in crowded scenes. *The Sci. World J.* **2014**, 12 pages (2014)
17. Hugg, J., Rafalin, E., Seyboth, K., Souvaine, D.: An experimental study of old and new depth measures. In: *Workshop on Algorithm Engineering and Experiments (ALENEX)*, pp. 51–64 (2006)
18. Icons. Icons Project (2000). <http://www.dcs.qmul.ac.uk/research/vision/projects/ICONS/>
19. Jiang, F., Wu, Y., Katsaggelos, A.K.: Detecting contextual anomalies of crowd motion in surveillance video. In: *16th IEEE International Conference on Image Processing*, pp. 1117–1120 (2009)
20. Liu, R.: On a notion of data depth based on random simplices. *The Ann. Stat.* **18**(1), 405–414 (1990)
21. Loménie, N., Stamon, G.: Morphological mesh filtering and α -objects. *Pattern Recogn. Lett.* **29**(10), 1571–1579 (2008)
22. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vis.* **60**(2), 91–110 (2004)
23. Mehran, R., Oyama, A., Shah, M.: Abnormal crowd behavior detection using social force model. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 935–942 (2009)
24. Miyajima, K., Ralescu, A.: Spatial organization in 2D images. In: *Third IEEE Conference on Fuzzy Systems*, pp. 100–105 (1994)
25. Pele, O., Werman, M.: The quadratic-chi histogram distance family. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) *ECCV 2010, Part II. LNCS*, vol. 6312, pp. 749–762. Springer, Heidelberg (2010)
26. PETS (2006). <http://www.cvg.rdg.ac.uk/PETS2006/data.html>
27. PETS (2009). <http://www.cvg.rdg.ac.uk/PETS2009/a.html>
28. Piciarelli, C., Micheloni, C., Foresti, G.L.: Trajectory-based anomalous event detection. *IEEE Trans. Circ. Syst. Video Technol.* **18**(11), 1544–1554 (2008)
29. Rubner, Y., Tomasi, C., Guibas, L.J.: The earth mover’s distance as a metric for image retrieval. *Int. J. Comput. Vis.* **40**(2), 99–121 (2000)
30. Saleemi, I., Shafique, K., Shah, M.: Probabilistic modeling of scene dynamics for applications in visual surveillance. *IEEE Trans. Pattern Anal. Mach. Intell.* **31**(8), 1472–1485 (2009)
31. Tissainayagam, P., Suter, D.: Object tracking in image sequences using point features. *Pattern Recogn.* **38**(1), 105–113 (2005)
32. Tran, D., Yuan, J., Forsyth, D.: Video event detection: From subvolume localization to spatio-temporal path search. *IEEE Trans. Pattern Anal. Mach. Intell.* **36**(2), 404–416 (2014)
33. Tukey, J.W.: Mathematics and the picturing of data. In: *International Congress of Mathematicians*, vol. 2, pp. 523–531 (1975)
34. Vardi, Y., Zhang, C.-H.: The multivariate l1-median and associated data depth. *Nat. Acad. Sci.* **97**(4), 1423–1426 (2000)
35. Visam. Visam Project (1997). <http://www.cs.cmu.edu/~vsam/>
36. Zhou, H., Yuan, Y., Shi, C.: Object tracking using SIFT features and mean shift. *Comput. Vis. Image Underst.* **113**(3), 345–352 (2009)