

Published in IET Computer Vision  
Received on 30th November 2007  
Revised on 14th May 2008  
doi: 10.1049/iet-cvi:20070062

In Special Issue on Visual Information Engineering



ISSN 1751-9632

# Extraction of activity patterns on large video recordings

L. Patino<sup>1</sup> H. Benhadda<sup>2</sup> E. Corvee<sup>1</sup> F. Bremond<sup>1</sup>  
M. Thonnat<sup>1</sup>

<sup>1</sup>INRIA, 2004 Route des Lucioles, Sophia Antipolis 06902, France

<sup>2</sup>Thales Communication, 160 Boulevard de Valmy, Colombes 92704, France

E-mail: jlpatino@sophia.inria.fr

**Abstract:** Extracting the hidden and useful knowledge embedded within video sequences and thereby discovering relations between the various elements to help an efficient decision-making process is a challenging task. The task of knowledge discovery and information analysis is possible because of recent advancements in object detection and tracking. The authors present how video information is processed with the ultimate aim to achieve knowledge discovery of people activity and also extract the relationship between the people and contextual objects in the scene. First, the object of interest and its semantic characteristics are derived in real-time. The semantic information related to the objects is represented in a suitable format for knowledge discovery. Next, two clustering processes are applied to derive the knowledge from the video data. Agglomerative hierarchical clustering is used to find the main trajectory patterns of people and relational analysis clustering is employed to extract the relationship between people, contextual objects and events. Finally, the authors evaluate the proposed activity extraction model using real video sequences from underground metro networks (CARETAKER) and a building hall (CAVIAR).

## 1 Introduction

Nowadays, more than ever, the technical and scientific progress requires human operators to handle large quantities of data. To treat this huge amount of records, the data-mining field can provide adequate solutions to synthesise, analyse and extract valuable information, which is generally hidden in the raw data. Applying data-mining techniques in large amounts of video data is now possible mainly because of the advance made in the field of object detection and tracking [1]. Data mining on video data has mainly been employed for annotation/retrieval processes [2–5]. The task consists in mining multiple visual features into categories associated with meaningful semantic keywords that will allow the retrieval of the video. Usually low level features such as colour, texture, shape and motion information are employed. A recent review on video retrieval can be found in [6]. The structured representation issued from the mining procedure gives a domain-dependent association that tries

to solve the well-known gap between low-level features and high-level concepts. Recently, particular attention has been turned to the trajectory information associated with mobile objects observed in the video. It is because, on the one hand, trajectory descriptors have been shown to be very useful on their own for video indexing and retrieval [7], and on the other hand the application of data-mining and machine-learning techniques to the study of trajectories has started to show its importance for activity understanding. This kind of analysis comes as a complement to current video monitoring/surveillance systems such as PRISMATICA [8], VISOR-BASE [9] or ADVISOR [10], which were rather oriented towards the real-time recognition of events of interest (fighting between persons, vandalism, a person jumping above a barrier, a group of people blocking an exit, and overcrowding situations). Although these systems begin to recognise robustly predefined events in the video, data mining/knowledge discovery on the activities contained in the video has not been addressed.

In this paper we show the procedure to achieve knowledge discovery to obtain meaningful trajectory patterns from video-data. A hierarchical agglomerative clustering algorithm is employed for this purpose. Moreover, we show how semantic meaning can be obtained from these patterns. For instance we can study the dynamics of the people characterised by a given trajectory type. We have proposed an effective knowledge modelling format which allows us to combine the motion pattern of a person or a group of persons with other features describing their interactions with the environment. We also show in our contribution how we can further apply data-mining techniques on this information. Specifically, we employ the relational analysis [11] for this purpose and show how behavioural patterns of interaction can be extracted. The techniques employed have been evaluated on annotated videos from the CAVIAR project [12] and on large underground video recordings (GTT metro, Torino, Italy and ATAC metro, Roma, Italy) from the CARETAKER project ([www.ist-caretaker.org](http://www.ist-caretaker.org)). These videos are associated with manually generated ground-truth.

This research has been done in the framework of the CARETAKER project, which is a European initiative to provide an efficient tool for the management of large multimedia collections. Such systems could be used in applications such as surveillance and safety issues, in urban planning, resource optimisation, elderly person monitoring. This work was partially presented in [13]. In this paper we employ an analysis on more features extracted from mobile objects detected in the videos. We present a new evaluation for our trajectory-clustering algorithm and we present results from both Torino and Roma underground sites.

The rest of the paper is structured as follows. In Section 2 we present the overall architecture of the proposed approach. While the video analysis system to track objects and detect events of interest in the video is explained in Section 3, the clustering of trajectories from tracked objects is detailed in Section 4. The proposed knowledge representation format is presented in Section 5. The relational analysis applied on the trajectory and the contextual information is explained in Section 6. Results are presented in Section 7. The proposed method is assessed in Section 8.

## 2 General structure of the proposed approach

The monitoring system is mainly composed of two different processing components (shown in Fig. 1). The first one is an on-line analysis subsystem for the real-time detection of objects and events previously defined in an ontology. This is a processing that goes on a frame-by-frame basis. At this level, detected events already contain semantic information describing people behaviour and interactions with the contextual objects of the scene. The second subsystem works off-line and achieves the extraction of activity

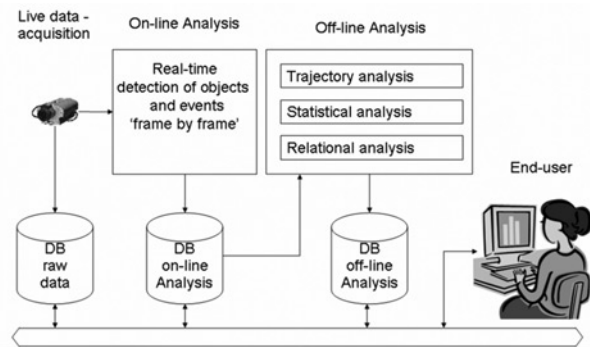


Figure 1 General architecture

patterns from the video. This subsystem is composed of three modules: the trajectory analysis module where we perform the clustering of trajectories, the object statistical analysis module where we compute meaningful measures on the object dynamics and the relational analysis module where we obtain behavioural patterns of interaction. For the storage of video streams and the metadata obtained after both on-line video processing and off-line analysis, three different databases (DB) exist: raw database (audio/video streams), on-line analysis database (tracked objects and real-time detected events), off-line analysis database (data mining-based events). In this work we only consider the video analysis although audio data has also been acquired in the project.

Streams of video are acquired at a speed of 5–25 frames/s, then directly stored in the raw database. The on-line analysis subsystem takes its input directly from the data acquisition component. Objects and events of interest are detected in real time, and tracking results are written in the on-line database at a speed of 5 frames/s. The on-line system triggers alarms for the security operators to take immediate actions.

The off-line analysis subsystem takes its input from the on-line database. This subsystem is dedicated to the manager or designer who wants to obtain global and long-term information from the monitored site. The user can specify a period of time where he/she wishes to retrieve and analyse stored information. In particular the user can access all databases to visualise specific events, streams of video and off-line information.

## 3 Real-time object/event detection

The first task of our data-mining approach is to detect in real time objects and events of interest. A brief description on how both objects and events of interest are detected in this work is given next. It must be noted that our data-mining approach can however be applied with any other detection and tracking technique.

### 3.1 Object detection

Assuming objects in a scene are individually detected in the image plane; tracking these objects would be a straightforward task to perform and reliable trajectories would be obtained. However, the efficiency of a tracking algorithm directly depends on the quality of the detected objects which is very sensitive to many factors such as the quality of the image (level of compression, accumulated noise throughout the optical device), dynamic occlusions (when several mobile projections onto the image plane overlap) [14], and the complexity of interaction of the objects evolving in a scene as we describe below.

Detecting objects in an image is a difficult and challenging task. Several algorithms have been proposed; two recent surveys [1, 15] include the research from the last 10 years in this regard. One solution widely employed consists of performing a thresholding operation between the pixel intensities of this frame with the pixel intensities of the background frame. The background image can be a captured image of the same scene having no foreground objects, or no moving objects in front of the camera.

The result of the thresholding operation is a binary mask of foreground pixels. The neighbouring foreground pixels are grouped together to form regions often referred to as 'blobs' which correspond to the moving regions in the image. If the 'moving objects' projection in the image plane do not overlap with each other, i.e. no 'dynamic occlusion', then each detected moving blob corresponds to a single moving object. However, as soon as occlusion occurs between objects, their moving regions fuse and separating them is not an easy task to do. Adding more informational cues about the detected objects, such as their colour content, shape information or multiple view tracking, increases the likelihood in separating the occluded regions. The detailed description of the background subtraction algorithm, which also estimates when the background reference image needs to be updated, can be found in [16].

Having 3D information about the scene under view enables the calibration of the camera. Point correspondences between selected 3D points in the scene and their corresponding point in the 2D image plane allow us to generate the 3D location of any points belonging to moving objects. Thus, the 3D data (i.e. width and height) of each detected moving blob can be measured as well as their 3D location on the ground plane in the scene with respect to a chosen coordinate system. The 3D object information is then compared against several 3D models defined by the user. From this comparison, a detected object is linked to a semantic class. For example, we have chosen a human being to be of average size: 170 cm in height and 60 cm in width. Smaller sizes are chosen for luggage items, and bigger sizes for groups of persons and very large sizes for crowds. The description and use of these kinds of 3D models can be found in [17]. The noisy

detected objects, associated with noisy 3D models, are classified according to the closest model they belong to.

### 3.2 Object tracking

Detected and classified 3D objects evolving in a scene can be tracked within the scope of the camera using the 3D information of their locations on the ground as well as their 3D dimensions. Tracking a few objects in a scene can be easy as far as they do not interact heavily in front of the camera: i.e. occlusion is rare and short. However, the complexity of tracking several mobile objects becomes a non-trivial and very difficult task to achieve when several objects' projected images overlap with each other on the image plane. Occluded objects have missing 3D locations, which create incoherency in the temporal evolution of their 3D locations.

Our tracking algorithm [18, 19] builds a temporal graph of connected objects over time to cope with the problems encountered during tracking. The detected objects are connected between each pair of successive frames by a frame-to-frame (F2F) tracker. Links between objects are associated with a weight (i.e. a matching likelihood) computed from three criteria: the similitude between their semantic classes, 3D dimensions and differences in 3D distance in the ground plane.

The graph of linked objects provided by the F2F tracker is then analysed by the tracking algorithm, also referred to as the long-term tracker, which builds paths of each mobile object according to the link features. The best path is then taken out as the trajectory of the related mobile objects.

### 3.3 Event detection

Events of interest for the user are created by the user himself according to a specific semantic language introduced by Vu *et al.* [20]. This language allows the user to use a designed ontology to detect from simple to complex events. The ontology is the set of all concepts relative to video events and of all the relations between concepts. There are two main types of concepts to be represented: physical objects (including physical objects of interest to be observed in a scene, also called as 'mobile objects' and 'contextual objects', which are defined by the user) and video events occurring in this scene related to the objects of interest.

A physical object of interest ' $\sigma$ ' is a physical object evolving in the scene, whose semantic class (i.e. person, group, crowd and luggage) is predefined by end-users and whose motion cannot be foreseen using a priori information. The tracked object is characterised by 2D and 3D features (e.g. a 3D location, width and height), a trajectory and an identifier. Using the 2D and 3D features, the object classification algorithm compares the object attributes with the predefined semantic classes (i.e. person, group, crowd and luggage) and assigns the corresponding semantic label to

the tracked object. A contextual object is a physical object attached to the scene, usually an equipment 'eq', or a zone of interest 'z'. The contextual object is usually static and whenever in motion, its motion can be foreseen using *a priori* information. For instance, the movements induced by a door or a chair can be foreseen.

A video event describes any event, action or activity happening in the scene and visually observable by cameras. Video events are characterised by the involved objects of interest (as described above), and their starting–ending times. Examples of events are 'detection of a person inside a zone', 'detection of an abandoned bag' and 'a meeting between two people'. For instance 'abandoned bag' consists in the detection of a bag with nobody around for a certain time. For event detection we have ourselves taken inspiration from methods published in PETS2006 workshop [21], and the specific implementation we have in our system can be found in [10].

We distinguish four types of video events, i.e. 'primitive state', 'composite state', 'primitive event' and 'composite event' which are classified into two categories, i.e. 'state' and 'event' defined as follows:

- A 'state' is a spatio-temporal property of a physical object valid at a given instant or constant on a time interval. A state characterises one or several physical objects of interest (e.g. person, crowd or vehicle) with or without respect to other physical objects.
- A 'primitive state' is a state which is directly inferred from visual attributes of physical objects computed by perceptual components. Usually, visual attributes have a numerical value and can correspond to general physical object properties for most video understanding applications. For example: 'A is inside a zone'.
- A 'composite state' is a combination of states. We call 'state components' all the sub-states composing the state and we call 'constraints' all the relations involving its components and its physical objects. For example: 'Person  $p_1$  is close to machine  $m$  and person  $p_2$  stays inside zone  $z$ '.
- An 'event' is one or several change(s) of state values at two successive time instants or on a time interval.
- A 'primitive event' is a change of primitive state values. Primitive events are more abstract than states but they represent the finest granularity of events. For example: 'Person  $p$  moves from zone  $z_1$  to zone  $z_2$ '.
- A 'composite event' is a combination of states and events. Usually, most abstract composite events have a symbolical/Boolean value and are directly linked to the goals of the given application. We call 'event components' all the sub-states/events composing the event and we call 'constraints'

all the relations involving its components and its physical objects.

In the applications of the work presented in this paper, the following events have been defined:

- $\text{inside\_zone}(o, z)$ : when an object ' $o$ ' is in the zone ' $z$ '.
- $\text{'stays\_inside\_zone}(o, z, T_1)$ ': when the event ' $\text{inside\_zone}(o, z)$ ' is being detected successively for at least  $T_1$  seconds
- $\text{'close\_to}(o, \text{eq}, D)$ ': when the 3D distance of an object location on the ground plane is less than the maximum distance allowed,  $D$ , from an equipment object 'eq'
- $\text{'stays\_at}(o, \text{eq}, D, T_2)$ ': when the event ' $\text{close\_to}(o, \text{eq}, D)$ ' is being consecutively detected for at least  $T_2$  seconds.
- $\text{'crowding\_in\_zone}(\text{crowd}, z)$ ': when the event ' $\text{stays\_inside\_zone}(\text{crowd}, z, T_3)$ ' is detected for at least  $T_3$  seconds.

Also, we have employed the following variables:

For mobile objects:

- object  $o = \{p, g, c, l, t, u\}$  with  $p$  = person,  $g$  = group,  $c$  = crowd,  $l$  = luggage,  $t$  = train and  $u$  = unknown.

For contextual objects:

- zone  $z = \{\text{platform}, \text{validating\_zone}, \text{vending\_zone}\}$
- equipment  $\text{eq} = \{g_1, \dots, g_{10}, \text{vm}_1, \text{vm}_2\}$  where  $g_i$  is the  $i$ th gate and  $\text{vm}_i$  is the  $i$ th vending machine.

Event thresholds:

- $T_1 = 60$  s,  $D = 1.50$  m,  $T_2 = 5$  s,  $T_3 = 120$  s.

$D$  corresponds to the Euclidean distance between 3D points of people position, given by the contact point of the person with the ground floor, and the 3D equipment localisation.

## 4 Trajectory analysis

The second layer of analysis in our approach is related to the knowledge discovery of higher semantic events from off-line analysis of activity recorded over a period of time that can span, for instance, from minutes to a whole day. Patterns of activity are first extracted from the analysis of trajectories. Then, the knowledge representation format we propose coupled with the statistical analysis provides a rich overview of the activities in the scene. For the analysis of more complex relationships between the objects observed in the scene, we employ the relational analysis clustering technique.

#### 4.1 Trajectory analysis: background and related work

Data mining/knowledge discovery techniques applied to trajectory data extract patterns hidden on the raw video that are critical to find out relevant information about the motion behaviour of a person (or set of persons) and their interactions with contextual objects of the scene in the video. In this regard, probably the most active research area has been normal/abnormal behaviour detection. Piciarelli *et al.* [22] employ a splitting algorithm applied on very structured scenes (such as roads) represented as a zone hierarchy. The drawback of this approach is that it is difficult to generalise on other domains where trajectories have less structure inherited by the scene. Anjum and Cavallaro [23] employ PCA to reduce the dimensionality of trajectories. They analyse the PCA first two components of each trajectory together with their associated average velocity vector. Mean-shift, with these features, is employed to seek the local modes and generate the clusters. The modes associated with very few data points are considered as outliers. The outlier condition is set as the 5% of the maximum peak in the dataset, but again the drawback of the approach is that the analysis is adapted to highly structured scenes. Similarly, Naftel and Khalid [24] first reduce the dimensionality of the trajectory data employing discrete Fourier transform coefficients and then apply a self-organising map (SOM) clustering algorithm to find normal behaviour. Antonini and Thiran [25] transform the trajectory data employing independent component analysis, while the final clusters are found employing an agglomerative hierarchical algorithm. In these approaches it is however delicate to select the number of coefficients that will represent the data after dimensionality reduction. Calderara *et al.* [26] employ a k-medoids clustering algorithm on a transformed space modelling different possible trajectory directions to find groups of normal behaviour. Abnormal behaviour is detected as a trajectory that does not fit into the established groups, however the approach is validated with acted abnormal trajectory.

Data mining of trajectories has also been applied with statistical methods. Gaffney and Smyth [27] employed mixtures of regression models to cluster hand movements, although the trajectories were constrained to have the same length. Hidden Markov models (HMMs) have also been employed [28–30]. However, it has been observed that the structures and probability distributions of this kind of approach are highly domain dependent and require a tedious stage of parameter tuning [23].

All these techniques are interesting, but little has been said about the adequacy of the trajectory clusters and end user expectations.

More than trajectory clusters, in this paper we are interested with extracting meaningful activity clusters, which differ also from normal abnormal behaviour

extraction and where clustering techniques have also been employed, not only on trajectory data but also on event data [31–34]. Thus we show first how it is possible to obtain meaningful trajectory patterns from video-data by selecting a large set of features from the trajectory and then employing an agglomerative clustering algorithm. Second, this trajectory-clustering stage is coupled with statistical analysis to infer meaningful activities occurring in the scene. Moreover, we complement the trajectory analysis with relational analysis to find out complex activity patterns corresponding to higher semantic relations between variables.

#### 4.2 Trajectory analysis: proposed method

For the trajectory pattern characterisation of the object, we have selected a comprehensive, compact and flexible representation. It is suitable also for further analysis as opposed to many video systems. They actually store the sequence of object locations for each frame of the video, which is a cumbersome representation with no semantic information.

If the dataset is made up of  $N$  objects, the trajectory for object  $j$  in this dataset is defined as the set of points  $[x_j(t), y_j(t)]$  corresponding to the points with main direction changes;  $x$  and  $y$  are time series vectors whose length is not equal for all objects as the time they spend in the scene is variable. Two key points defining these time series are the beginning and the end,  $[x_j(1), y_j(1)]$  and  $[x_j(\text{end}), y_j(\text{end})]$  as they define where the object is coming from and where it is going to. We build a feature vector from these two points. Additionally, we also include the directional information given as  $[\cos(\theta), \sin(\theta)]$ , where  $\theta$  is the angle which defines the vector joining  $[x_j(1), y_j(1)]$  and  $[x_j(\text{end}), y_j(\text{end})]$ .

We feed the feature vector formed by these elements to a hierarchical clustering algorithm. For a data set made of  $N$  trajectories there are  $N \times (N - 1)/2$  pairs in the dataset. We employ the Euclidean distance as a measure of similarity to calculate the distance between all trajectory features. To avoid one feature to prime over the others, particularly because distances in  $x$  and  $y$  are much bigger than distances on direction, input data related to the beginning and end of a trajectory are first normalised. The mean for each feature in  $x$  and  $y$  is first removed, then each record has its value divided by the standard deviation calculated from each feature. Because  $[\cos(\theta), \sin(\theta)]$  are already bounded with values  $[-1, 1]$ , these directional features are not transformed. Other distances have been envisaged such as the weighted sum of the features. Object trajectories with the minimum distance are clustered together. When two or more trajectories are set together the mean of the trajectory features is taken into account for further clustering. The successive merging of clusters is listed by a dendrogram. The evaluation of the dendrogram is typically subjective by adjudging which distance threshold appears to create the most natural grouping of the data.

The final number of clusters is set by maximising the evaluation criteria, which is defined in the next section. As the acquisition is performed in a multi-camera environment the clusters obtained can be generalised to different camera views with a 3D calibration matrix stage carried out for the on-line analysis system.

### 4.3 Trajectory analysis: evaluation

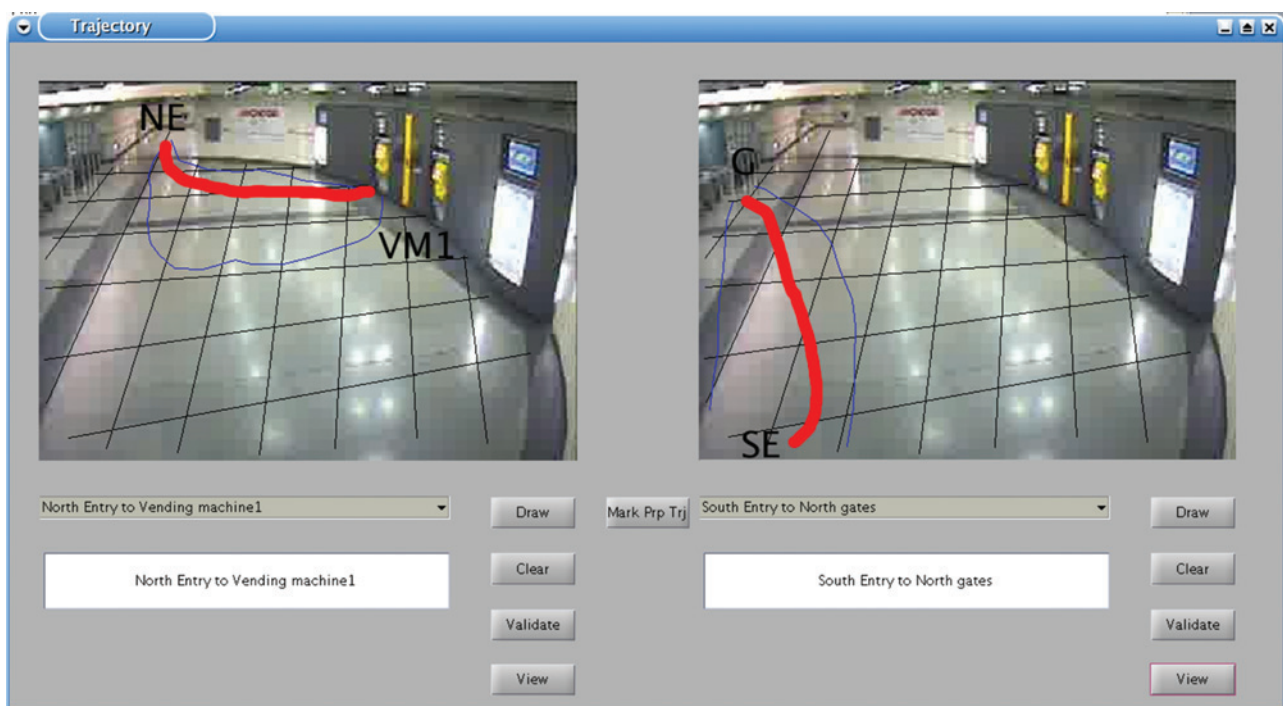
In order to evaluate our trajectory analysis approach, we have defined a Ground-truth data set containing over 300 trajectories. The ground-truth trajectories were manually drawn on an image illustrating the empty scene. Fig. 2 shows the empty scene for the Torino metro with some drawn trajectories. Semantic descriptions such as 'from north entry to vending machine1' were generated. There are 100 such annotated semantic descriptions, which are called trajectory types in the following. Each trajectory type is associated with a main trajectory that best matches that description. Besides, two complementary trajectories define the confidence limits within which we can still associate that semantic description. In Fig. 2 the main trajectory of each trajectory type is represented by a thick line while thin lines represent the complementary trajectories. Thus, each trajectory type is associated with triplets of trajectories.

We compute two performance measures to validate the quality of the proposed clustering approach, namely, confusion and dispersion. The former gives an indication of how many trajectory types of the ground-truth (how many main trajectories) are merged together in a single cluster

resulting from the agglomerative procedure. Ideally the clustering algorithm should be able to dissociate all main trajectories to separate all trajectory types. In this case we would have confusion = 1 (only one trajectory type is included in one cluster). If two main trajectories are included in the same cluster, then we say that confusion = 2 as two different trajectory types are merged in the same cluster. In general terms the confusion value for a cluster equals the number of main trajectories included.

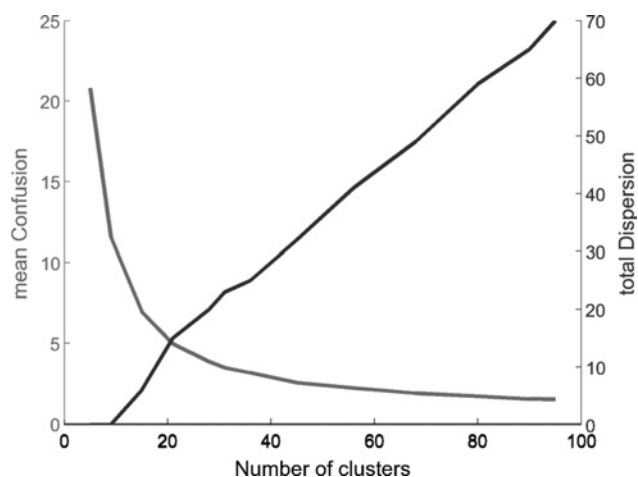
The latter performance measure (dispersion) indicates the number of erroneous clustered trajectories. This refers to the case where the main trajectory and any of its complementary trajectories (trajectories defining the confidence limits of the main trajectory) have been splitted into different clusters. Each 'left-apart' complementary trajectory increments the Dispersion measure by one unit. Fig. 3 depicts the evolution of these two factors depending on the number of clusters which is chosen when running the clustering algorithm. For instance, for 21 clusters we have in total 15 complementary trajectories badly clustered (dispersion) and 5 main trajectories per cluster (mean confusion)

Because all trajectories cannot be equally observed by the camera (for instance distinguishing all turnstiles in the upper left corner would require a larger spatial resolution), it is actually very difficult to achieve a bijection between the semantic labels and the resulting clusters. However, we aim at having the lowest possible confusion level together with the lowest percentage of dispersion. From Fig. 3, it can be



**Figure 2** Ground-truth for two different semantic clusters

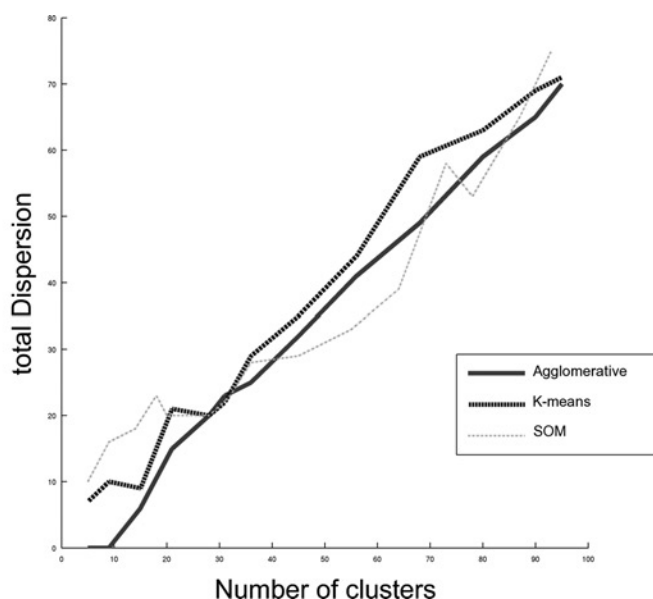
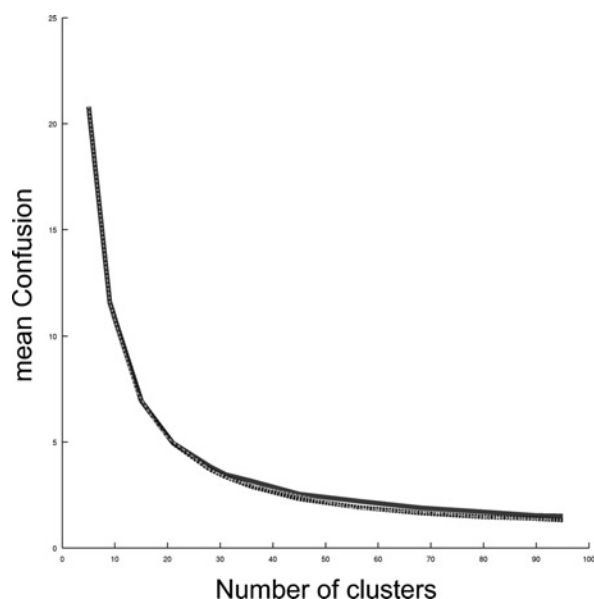
Left panel shows trajectories associated with the trajectory type 'From North Entry (NE) to Vending Machine1 (VM1)'  
Right panel shows trajectories associated with the trajectory type 'From South Entry (SE) to Gates (G)'



**Figure 3** Evolution of the clustering quality measures Confusion (fusion of ground-truth semantic labels) and Dispersion (erroneous clustered trajectories) as a function of the number of clusters

observed that a good compromise is achieved for a number of clusters between 15 and 35.

We also compared the agglomerative clustering algorithm in our application to other well-known clustering techniques such as  $k$ -means [35] and Kohonen SOM [36] employing the same ground-truth data. We observed that in the range of 1–35 clusters the Agglomerative algorithm has the less dispersion, whereas the confusion remains almost identical for the three techniques. When the number of clusters increases, the SOM algorithm has a smaller dispersion than the other two techniques. However, the amount of erroneously clustered data goes beyond 10%. The choice of the agglomerative algorithm remains thus valid for the



**Figure 4** Evolution of the clustering quality measures Confusion (fusion of ground-truth semantic labels) and Dispersion (erroneous clustered trajectories) as a function of the number of clusters for the agglomerative clustering algorithm,  $k$ -means and Kohonen SOM

number of clusters chosen before (between 15 and 35). Fig. 4 shows the evolution of the confusion and dispersion measures according to the number of clusters.

We further evaluated the three techniques with typical clustering validity indexes such as Silhouette, Dunn and Davis–Bouldin indexes, which are described next:

**4.3.1 Silhouette index:** The Silhouette index [37, 38] is defined as follows. Consider a data object  $v_j$ ,  $j \in \{1, 2, \dots, N\}$  belonging to cluster  $cl_i$ ,  $i \in \{1, \dots, c\}$ . This means that object  $v_j$  is closer to the prototype of cluster  $cl_i$  than to any other prototype. Let the average distance of this object to all objects belonging to cluster  $cl_i$  be denoted by  $a_{ij}$ . Also, let the average distance of this object to all objects belonging to another cluster  $i'$ ,  $i' \neq i$ , be called  $d_{i'j}$ . Finally, let  $b_{ij}$  be the minimum  $d_{i'j}$  computed over  $i' = 1, \dots, c$ , which represents the dissimilarity of object  $j$  to its closest neighbouring cluster. The Silhouette index is then

$$S = \frac{1}{N} \sum_{j=1}^N s_j, \quad \text{where } s_j = \frac{b_{ij} - a_{ij}}{\max\{a_{ij}, b_{ij}\}}$$

This way, the best partition is achieved when  $S$  is maximised, which implies minimising the intra-cluster distance ( $a_{ij}$ ) while maximising the inter-cluster distance ( $b_{ij}$ ).

**4.3.2 Dunn index:** The Dunn index [39] is defined as follows. Let  $cl_i$  and  $cl_{i'}$  be two different clusters of the input dataset. Then, the diameter  $\Delta$  of  $cl_i$  is defined as

$$\Delta(cl_i) = \max_{v_j, v_{j'} \in cl_i} \{d(v_j, v_{j'})\}$$

Let  $\delta$  be the distance between  $cl_i$  and  $cl_{i'}$ . Then  $\delta$  is defined as

$$\delta(cl_i, cl_{i'}) = \min_{v_j \in cl_i, v_{j'} \in cl_{i'}} \{d(v_j, v_{j'})\}$$

and  $d(x, y)$  indicates the distance between points  $x$  and  $y$ .

For any partition, the Dunn index is

$$v_D = \frac{\min_i \left\{ \min_{i'} \left\{ \frac{\delta(cl_i, cl_{i'})}{\max_i \{cl_i\}} \right\} \right\}}{\max_i \{cl_i\}} \text{ and } i, i' \in \{1, \dots, N\}, i' \neq i$$

Larger values of  $v_D$  correspond to a good clustering partition.

**4.3.3 Davis–Bouldin index:** The Davis–Bouldin index [40] is defined as follows: This index is a function of the ratio of the sum of within-cluster scatter to between-cluster separation. The scatter within cluster,  $cl_i$ , is computed as

$$S_i = \frac{1}{|cl_i|} \sum_{v_j \in cl_i} \{\|v_j - m_i\|\}$$

$m_i$  is the prototype for cluster  $cl_i$ . The distance  $\delta$  between clusters  $cl_i$  and  $cl_{i'}$  is defined as

$$\delta(cl_i, cl_{i'}) = \|m_i - m_{i'}\|$$

The Davies–Bouldin (DB) index is then defined as

$$DB = \frac{1}{N} \sum_{i=1}^N R_i \quad \text{with } R_i = \max_{i'} R_{ii'} \quad i, i' \in \{1, \dots, N\} \\ i' \neq i$$

and

$$R_{ii'} = \frac{S_i + S_{i'}}{\delta(cl_i, cl_{i'})}$$

Low values of the DB index are associated with a proper clustering.

The results from applying these measures are summarised in Table 1 for a number of clusters equal to 31, chosen as best

choice from the confusion dispersion curves. It can be observed that the agglomerative algorithm is the best choice according to Silhouette and Dunn indexes. According to DB index, the SOM technique would be the best choice. However, for this index and the Silhouette, all methods are close from each other. For the Dunn index, the agglomerative method is clearly the best choice.

## 5 Knowledge representation and statistical analysis

There are two main types of concepts to be represented from the video: physical objects of the observed scene and video events occurring in the scene. The former category can still be further subdivided into two types of physical objects of interest: mobile and contextual objects. Mobile objects of interest are the source of action occurring in the scene. Contextual objects are 3D objects of the empty scene model corresponding to the static environment of the scene. For the off-line analysis of both types of concepts, and with the aim of setting the data in a suitable format to achieve knowledge discovery, we separate the information corresponding to the activities occurring over a period of time on three different semantic tables, namely mobile objects, contextual objects and video events. This information is characterised by a set of specific features, which is currently being enriched along the CARETAKER project. We are reporting in the following the features that have been used in our experimentations.

### 5.1 Mobile objects

The mobile object,  $m$ , can be represented as a ten-tuple  $m = \{m^{\text{id}}, m^{\text{type}}, m^{\text{start}}, m^{\text{end}}, m^{\text{duration}}, m^{\text{dist\_org\_dest}}, m^{\text{shape}}, m^{\text{involved\_events}}, m^{\text{significant\_event}}, m^{\text{trajectory}}\}$

where  $m_j$  with  $j \in \{1, 2, \dots, N\}$  is then

- $m_j^{\text{id}}$ : the identifier label of the object.  $m_j^{\text{id}} \in \mathbb{Z}^+$
- $m_j^{\text{type}}$ : the class the object belongs to:
- $m_j^{\text{type}} \in \{\text{'Person'}, \text{'Group'}, \text{'Crowd'}, \text{'Train'}, \text{'Luggage'}\}$ .
- $m_j^{\text{start}}$ : the time when the object is first seen.

**Table 1** Evaluation of the clustering quality indexes (Silhouette, Dunn and Davis–Bouldin) for the agglomerative clustering algorithm, k-means and Kohonen SOM

|                |               | Clustering technique |          |         |                   |
|----------------|---------------|----------------------|----------|---------|-------------------|
|                |               | Agglomerative        | K-means  | SOM     |                   |
| validity index | Silhouette    | 0.32116              | 0.29501  | 0.3036  | higher is better  |
|                | Dunn          | 0.10098              | 0.064249 | 0.05031 | higher is better  |
|                | Davis–Bouldin | 0.58685              | 0.49352  | 0.37311 | smaller is better |



- $m_j^{\text{end}}$ : the time when the object is last seen.
- $m_j^{\text{duration}}$ : the total time the object has been observed.  
 $m_j^{\text{duration}} = m_j^{\text{end}} - m_j^{\text{start}}$
- $m_j^{\text{dist\_org\_dest}}$ : the total distance walked from origin to destination.
- $m_j^{\text{dist\_org\_dest}} = \sum_{t=1}^{T_j-1} \sqrt{(x_j(t+1) - x_j(t))^2 + (y_j(t+1) - y_j(t))^2}$   
and  $t \in \{1, \dots, T_j\}$  where  $T_j$  is the number of trajectory sampled points for the object  $m_j$ .
- $m_j^{\text{shape}}$ : the label describing the object's shape depending on the object's ratio height/width.  
 $m_j^{\text{shape}} = \{\text{'Small'}, \text{'Medium'}, \text{'Large'}\}$
- $m_j^{\text{involved\_events}}$ : all occurring Events related to the identified object.
- $m_j^{\text{significant\_event}}$ : the most significant event among all events. This is calculated as the most frequent event related to the mobile object.
- $m_j^{\text{trajectory}}$ : the trajectory cluster identifier characterising the object.  $m_j^{\text{trajectory}} = i \Leftrightarrow m_j \in \text{cl}_i$

### 5.2 Contextual objects

The activities involving a contextual object,  $c$ , can be characterised by a 12-tuple

$$c = \left\{ \begin{array}{l} c^{\text{id}}, c^{\text{type}}, c^{\text{start}}, c^{\text{end}}, c^{\text{involved\_events}}, c^{\text{significant\_event}}, \\ c^{\text{rare\_event}}, c^{\text{event\_histogram}}, c^{\text{involved\_mobiles}}, \\ c^{\text{mobiles\_histogram}}, c^{\text{use\_duration}}, c^{\text{mean\_time\_of\_use}} \end{array} \right\}$$

where  $c_k$  and  $k = \{1, \dots, O\}$  is then,

$c_k^{\text{id}}, c_k^{\text{involved\_events}}, c_k^{\text{significant\_event}}$  are defined in the same way as for the mobile objects but referring to contextual objects and  $c_k^{\text{type}} \in z \cup \text{eq}$ , as defined above  $z = \{\text{'platform'}, \text{'validating\_zone'}, \text{'vending\_zone'}\}$ ,  $\text{eq} = \{g_1, \dots, g_{10}, \text{vm}_1, \text{vm}_2\}$  where  $g_i$  is the  $i$ th gate and  $\text{vm}_i$  is the  $i$ th vending machine.

$c_k^{\text{start}}, c_k^{\text{end}}, c_k^{\text{involved\_mobiles}}$  are defined in the same way as for the mobile objects interacting with the contextual object.

The remaining fields indicate

- $c_k^{\text{rare\_event}}$ : this is the rarest event.
- $c_k^{\text{event\_histogram}}$ : gives the number of occurrence of all involved events per type of event.
- $c_k^{\text{mobiles\_histogram}}$ : gives the number of appearance for all involved mobile objects per object type.

- $c_k^{\text{use\_duration}}$ : percentage of occupancy (or use of a contextual object). For instance, the ticket machine has a 10% of use over the observation time.
- $c_k^{\text{mean\_time\_of\_use}}$ : average time for a mobile object to interact with the contextual object.

The contextual objects to be monitored are predefined by the end-users in the model of the scene environment. This modelling phase is a quick process and enables to acquire the end-user expertise on the objects of interest. For the video sequences analysed in this work, the contextual objects of interest are: 'platform hall', 'gates' and 'vending machines'.

### 5.3 Video events

A video event,  $e$ , can be represented as a 6-tuple

$$e = \{e^{\text{id}}, e^{\text{type}}, e^{\text{start}}, e^{\text{end}}, e^{\text{involved\_mobiles}}, e^{\text{involved\_context\_obj}}\}$$

where  $e_l$  and  $l = \{1, \dots, M\}$  is then,

- $e_l^{\text{id}}$ : the identifier label for the detected event.  $e_l^{\text{id}} \in \mathbb{Z}^+$
- $e_l^{\text{type}}$ : the class where the event belongs to.  $e_l^{\text{type}} = \{\text{'close\_to'}, \text{'stays\_at'}, \dots\}$
- $e_l^{\text{start}}$ : the first time when the event is detected.
- $e_l^{\text{end}}$ : the last time when the event is seen.
- $e_l^{\text{involved\_mobiles}}$ : the identifier label of the mobile objects involved in that event.
- $e_l^{\text{involved\_context\_obj}}$ : the contextual object involved in that event.

### 5.4 Statistical analysis

Statistical information can be obtained from the mobile objects and the contextual objects as well as their interactions. This is a major information source for the end-user. For instance, on large metro video recordings, there is spatial and temporal information on the use of contextual objects. In this work we calculate the following statistical indexes.

- **number\_of\_users**: the total number of people interacting with a contextual object. The number\_of\_users is calculated as follows.

Let  $E = \{e_{l'}\} | e_{l'}^{\text{involved\_context\_obj}} = c_k^{\text{type}}, \text{ and } l' \subseteq l. l = \{1, \dots, M\}$ .  $M$  is the total number of events observed in the video, thus  $E \subseteq \{e_{l'}\}$ .

$\text{number.of.users} = \text{cardinal}(c_k^{\text{involved.mobiles}})$  with  $c_k^{\text{involved.mobiles}} = \{e_{l'}^{\text{involved.mobiles}}\}$ .  $c_k^{\text{involved.mobiles}}$  contains no repeated elements.

- `percentage_of_use`: the ratio between the total period of time a contextual object is in use to the total observation period

$\text{percentage\_of\_use} = \frac{e_{j''}^{\text{end}} - e_{j''(1)}^{\text{start}}}{\text{observation\_period}}$  if  $\text{cardinal}(E) = M'$  and  $M' \leq M$ .  $M'$  is the total number of events observed in the video where the contextual object  $c_k$  appears.

- `interaction_duration`: the mean time a user spends when interacting with a contextual object.

Let  $E' = \{e_{j''}\}_{e_{j''}^{\text{involved\_context\_obj}} = c_k^{\text{type}}, e_{j''}^{\text{involved\_mobiles}} \subseteq c_k^{\text{involved\_mobiles}}}$

$\text{Interaction\_duration} = \text{mean}(e_{j''(M'')}^{\text{end}} - e_{j''(1)}^{\text{start}})$  if  $\text{cardinal}(E') = M''$  and  $M'' \leq M'$ .  $M''$  is the total number of events observed in the video where one of the objects listed in  $c_k^{\text{involved\_mobiles}}$  interacts with the contextual object  $c_k$ .

The statistics can be visualised with an interactive user interface enabling us to study a given variable such as zone of interest, equipment etc.

## 6 Relational analysis

Once all statistical measures of the activities in the scene have been computed and the corresponding information is put into the proposed model format, we aim at discovering complex relationships that may exist between mobile objects themselves, and between mobile objects and contextual objects in the scene. For this task, the clustering methodology we decided to use was relational analysis and regularised similarity (RARES). This methodology gathers two different technologies: relational analysis theory and regularised similarity [41–43]. Relational analysis has been initiated and developed at the European Centre of Applied Mathematics (ECAM) at IBM France by Marcotorchino and Michaud [11, 44]. The principle of relational analysis consists in transforming the data usually represented as a  $N \times M$  rectangular matrix where  $N$  is the number of objects to be clustered and  $M$  is the number of variables measured on these objects to two new  $N \times M$  matrices representing, respectively, a global similarity and dissimilarity measures for each pair of objects. The relational analysis algorithm will compare, for any two objects, their similarity and their dissimilarity. If the similarity is greater than the dissimilarity, then these two objects will be put in the same cluster of the final obtained partition. The output of the relational analysis algorithm is a set of groups of objects, where inside each group the objects are more similar to each other than to the objects belonging to another group.

To define the global similarity matrix, relational analysis transforms each variable  $V^k$  ( $k = 1, 2, \dots, M$ ) to a  $N \times M$  matrix  $S^k$  where the term  $s_{ii'}^k$  is the similarity measure between the two objects  $i$  and  $i'$  w.r.t. variable  $V^k$ . A dissimilarity measure  $\bar{s}_{ii'}^k$  is then computed as the

complement to the maximum similarity measure possible between these two objects. As the similarity between two different objects is less or equal to their self-similarities: that is  $s_{ii'}^k \leq \min(s_{ii}^k, s_{i'i}^k)$ , the dissimilarity is  $\bar{s}_{ii'}^k = \min(s_{ii}^k, s_{i'i}^k) - s_{ii'}^k$ , then we define that the global similarity measure between objects  $i$  and  $i'$  over the  $M$  variables is  $s_{ii'} = \sum_{k=1}^M s_{ii'}^k$  and their global dissimilarity is also  $\bar{s}_{ii'} = \sum_{k=1}^M \bar{s}_{ii'}^k$ . To cluster a population of  $N$  objects described by  $M$  variables, the relational analysis theory is based on the maximisation of the Condorcet criterion  $C(X) = \sum_{i=1}^N \sum_{i'=1}^N (s_{ii'} x_{ii'} + \bar{s}_{ii'} \bar{x}_{ii'})$  where  $X$  is a binary  $N \times M$  matrix representing the partition to discover in the data. The general term  $x_{ii'}$  of matrix  $X$  is defined as follows

$$x_{ii'} = \begin{cases} 1 & \text{if } i \text{ and } i' \text{ are in the same cluster} \\ 0 & \text{otherwise} \end{cases}$$

and  $\bar{x}_{ii'} = 1 - x_{ii'}$

The mathematical formulation of the criterion to be maximised is

$\max(C(X))$  w.r.t.

$$\begin{cases} x_{ii} = 1 & \text{reflexivity} \\ x_{ii'} = x_{i'i} & \text{symmetry} \\ x_{ii'} + x_{i'i''} - x_{ii''} \leq 1 & \text{transitivity} \end{cases}$$

It seems evident that variables having only two modalities, for example, will tend to generate more similarities than those variables having bigger number of modalities. For example, it is less likely, for two persons chosen at random in Paris, to live in the same district (20 symbolic values) than to have the same gender (only 2 symbolic values). Regularised similarity, developed by Benhadda and Marcotorchino [42, 43], is a theory taking into account these internal structures, during the computation of the individuals' similarities; to rebalance the influences (too strong or too weak) induced, in an implicit way, by these structures. This will favour certain variables compared with the others and will create, thus, some biases. Regularised similarity did not bring noticeable clustering improvements in our analysis. Indeed, the number of symbolic values between variables does not change dramatically. Simple similarity is reported in this paper. Relational analysis is employed here for activity clustering as it is able to characterise heterogeneous data (including symbolic and numerical attributes) contrary to hierarchical clustering, which works with numerical-only data.

In the present work, the input of the system is then the whole set of detected mobile objects,  $m_j$   $j \in \{1, 2, \dots, N\}$ , with the features described in section 5.1, namely:  $m_j = \{m_j^{\text{type}}, m_j^{\text{duration}}, m_j^{\text{dist\_org\_dest}}, m_j^{\text{shape}}, m_j^{\text{significant\_event}}, m_j^{\text{trajectory}}\}$ , thus in our analysis  $M = 6$ . The output of the analysis is the

final partition, named  $\Omega$ , which indicates the elements  $m_j$  and  $m_{j'}$  that should be set together in the same cluster  $\Pi$

Several indicators are computed during the clustering process, to build up the final partition, named  $\Omega$  and to measure the quality of the obtained results. For a couple of objects  $i$  and  $i'$  belonging to the population  $P$  and two clusters  $\Pi$  and  $\Pi'$  belonging to  $\Omega$ , we define:

- The maximum similarity possible  $\Lambda_{ii'}$  between any two objects  $i$  and  $i'$  is:  $\Lambda_{ii'} = \text{Min}(c_{ii}, c_{i'i})$
- the link  $L_{ii'}$  between two objects by:  $L_{ii'} = c_{ii'} - \alpha \times \Lambda_{ii'}$  where  $\alpha$  is a parameter such that  $0 < \alpha \leq 0.5$
- the link  $L_{i\Pi}$  between object  $i$  and cluster  $\Pi$ , by:  $L_{i\Pi} = \sum_{i' \in \Pi} L_{ii'}$
- the agreement  $A_{\Pi\Pi'}$  between the two clusters, by:  $A_{\Pi\Pi'} = \sum_{i \in \Pi} \sum_{i' \in \Pi'} c_{ii'}$
- the disagreement  $\bar{A}_{\Pi\Pi'}$  between two clusters, by:  $\bar{A}_{\Pi\Pi'} = \sum_{i \in \Pi} \sum_{i' \in \Pi'} \bar{c}_{ii'}$
- the maximal own similarity  $\Lambda_{\Pi\Pi'}$  between two clusters, by:  $\Lambda_{\Pi\Pi'} = \sum_{i \in \Pi} \sum_{i' \in \Pi'} \Lambda_{ii'}$
- and the link  $L_{\Pi\Pi'}$  between two clusters, by:  $L_{\Pi\Pi'} = A_{\Pi\Pi'} - \alpha \times \Lambda_{\Pi\Pi'}$

The quality  $Q_{\Pi}$  of a particular cluster  $\Pi$  is defined by

$$Q_{\Pi} = \frac{A_{\Pi\Pi} + 2 \times \sum_{\Pi' \neq \Pi} \bar{A}_{\Pi\Pi'}}{\Lambda_{\Pi\Pi} + 2 \times \sum_{\Pi' \neq \Pi} \Lambda_{\Pi\Pi'}}$$

This measure takes into account at the same time the inner homogeneity and the external heterogeneity of cluster  $\Pi$ . The more the objects belonging to  $\Pi$  are similar to each other and in the same time the more they are dissimilar from the objects belonging to the other clusters, the better is the quality of the cluster.

The quality  $Q$  of the final partition  $\Omega$  is an indicator that measures the total coherence of  $\Omega$ , and is given by the formula

$$Q = \frac{\sum_{\Pi \in \Omega} (A_{\Pi\Pi} + \sum_{\Pi' \neq \Pi} \bar{A}_{\Pi\Pi'})}{\sum_{\Pi \in \Omega} \sum_{\Pi' \in \Omega} \Lambda_{\Pi\Pi'}}$$

## 7 Results

### 7.1 Results with annotated data

We first tested the validity of our clustering algorithms (hierarchical and relational clustering) on labelled video data. CAVIAR is an EC-funded project that has made

available a dataset of video clips with hand-labelled ground truth [12]. We focused our attention on the first part of the dataset, which contains people observed at the lobby entrance of a building. The annotated ground-truth include for each person its bounding box (id, centre coordinates, width, height, main axis orientation) with a description of his/her movement type (inactive, active, walking, running) for a given situation (moving, inactive, browsing) and with a given scenario context (browsing, immobile, left object, walking, drop down).

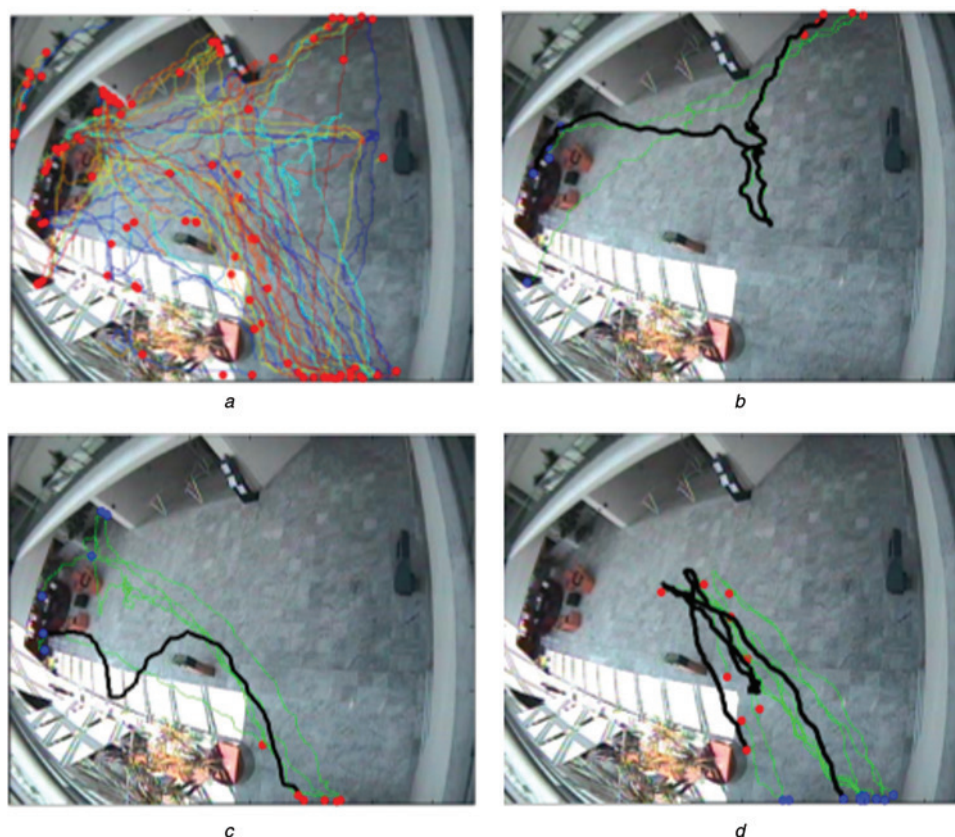
We have applied first the hierarchical clustering algorithm to a dataset containing 164 persons or objects. We have extracted their trajectories as the centre coordinates of the bounding box over time. We have tuned the algorithm with the user interface to obtain 31 meaningful clusters. Fig. 5a shows all trajectories from this dataset. The remaining plots in the figure are the three most common paths undertaken. As it can be observed, clear trajectory patterns can be extracted from the clusters. For instance, the three paths clusters from Fig. 4 can be labelled: cluster 8: 'entering right and up/exiting left' (Fig. 5b), cluster 11: 'entering right bottom/exiting left' (Fig. 5c) and cluster 15: 'entering right-middle/exiting right-bottom' (Fig. 5d).

We have then applied relational analysis to obtain higher relations between objects. For this purpose we have employed the object representation format described before. To perform the evaluation and because the annotated ground-truth was available with a situation and context description, we have generated a set of events by concatenating the three pieces of available information: movement's type ( $T$ ), context ( $C$ ) and situation ( $S$ ). The different symbolic values that the 3-tuple TCS can take are presented in Table 2.

For example, an event having the value 'awm' is related to an object with movement's type = '(a)ctive', within a context = '(w)alking' and situation = '(m)oving'. In the analysed data, an object is usually involved in between 1 and 12 such events during the observation time. To give account of the temporal succession of events, the 3-tuple TCS is sequentially numbered as the events appear. For instance if three events are in succession associated to a mobile object, these events will be designated as [T1 C1 S1], [T2 C2 S2] and [T3 C3 S3]. A portion of the input data is shown below in Table 3.

One of the clusters that RARES has discovered in the CAVIAR dataset is presented in Fig. 6. This cluster contains four detected objects. All objects are involved in only one event corresponding to inactive objects, in an inactive situation and in an immobile context. This cluster actually corresponds to the objects 'bag' that are annotated in the CAVIAR database as abandoned.

Another cluster is presented in Fig. 7. In this case all items are people that are involved in at least three events. For all



**Figure 5** Trajectories detected in the CAVIAR dataset and three clusters showing most common undertaken paths

- a Original set of CAVIAR trajectories
- b Cluster 8
- c Cluster 11
- d Cluster 15

Beginnings of the trajectories are indicated by red points  
Ends of trajectories are indicated by blue points

people, the first event is of walking type, in a moving situation and within a context where something drops\_down (75% of the cases). Then the situation and the context will evolve and all people will be involved in an event (third event), described with an inactive situation, within a context where something drops\_down and with an inactive movement type (75% of the cases). This cluster matches actually the objects that were annotated in the CAVIAR database as 'falls down' and included in the fighting (one man down) scenarios.

**Table 2** Semantic event information in CAVIAR

| Type       | Context       | Situation  |
|------------|---------------|------------|
| (i)nactive | (b)rowsing    | (m)oving   |
| (a)ctive   | (i)mmobile    | (i)nactive |
| (w)alking  | (l)eft object | (b)rowsing |
| (r)unning  | (w)alking     |            |
|            | (d)rop down   |            |

The activities described in the CAVIAR ground-truth are thus correctly retrieved with our algorithm. Table 4 gives a quantitative evaluation on the correspondence between the events annotated in the CAVIAR ground-truth and the clusters obtained by the relational analysis algorithm. As it can be observed, all instances from the left luggage event are well recognised in cluster 4. Fighting situations are generally recognised in cluster 0. The missing fighting cases are those corresponding to the case where one of the individuals involved in the fight falls down and lies on the floor. This kind of situation is matching cluster 5. Most cases of the browsing event are matched in cluster 2.

## 7.2 Results with large video recordings

We have processed 73 000 frames of video from the Torino Underground (GTT, Italy), with an acquisition rate of 25 frames/s equalling ~50 min of video. We have analysed off-line this period of time. We have applied the hierarchical clustering algorithm on the trajectories of mobile objects to obtain the common behavioural paths undertaken by the people in the hall. Fig. 8 presents the whole dataset of trajectories that we have analysed. Using our interactive user interface to maximise the evaluation

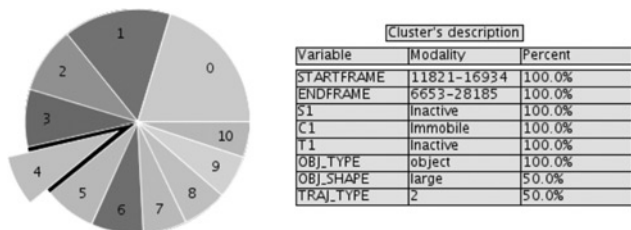
**Table 3** Input matrix used for the relational analysis clustering method

| Obj_id | Obj_type | Startframe | Endframe | Trajectory_type | Obj_shape | Involved_Event1 | Involved_Event2 | Involved_Event3 |
|--------|----------|------------|----------|-----------------|-----------|-----------------|-----------------|-----------------|
| 84     | Person   | 3785       | 14312    | 19              | big       | wwm             |                 |                 |
| 85     | Person   | 4338       | 14611    | 6               | small     | iii             | aii             | wirr            |
| 86     | Person   | 4459       | 14523    | 16              | tall      | wwm             |                 |                 |
| 88     | Person   | 4463       | 14684    | 7               | big       | rwm             | wwm             |                 |
| 87     | Person   | 4491       | 14684    | 10              | big       | wm              | aii             |                 |
| 89     | Person   | 4685       | 15097    | 15              | tall      | wm              | iii             | aii             |
| 90     | Person   | 4757       | 14938    | 16              | tall      | wwm             |                 |                 |
| 93     | Person   | 5054       | 15542    | 16              | big       | wm              | aii             | wirr            |
| 92     | Person   | 5255       | 15547    | 3               | small     | iii             |                 |                 |
| 94     | Person   | 5548       | 15670    | 15              | large     | wm              | aii             | iii             |
| 95     | Person   | 5695       | 15750    | 14              | small     | iii             | aii             | iii             |
| 96     | Person   | 5817       | 16699    | 14              | big       | abb             | wbm             | abb             |
| 98     | Person   | 6051       | 16563    | 19              | small     | iii             |                 |                 |
| 99     | Person   | 6920       | 17082    | 14              | small     | wm              | aii             | iii             |
| 103    | Person   | 6942       | 17513    | 4               | small     | wm              | aii             | iii             |
| 104    | Person   | 6966       | 17513    | 4               | small     | wm              | aii             | iii             |

Note that only some objects from the total detected set are represented in the table

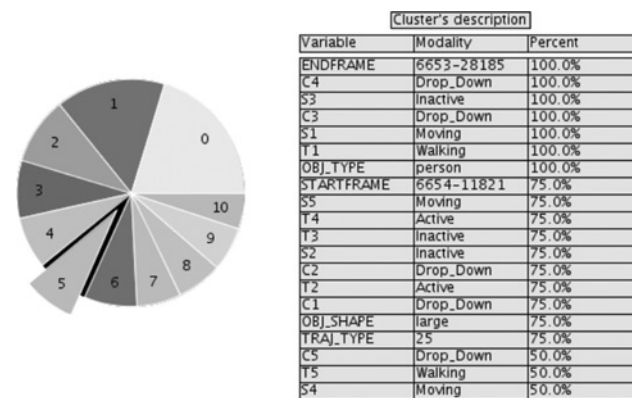
criteria (explained in Section 4.3), we have applied the hierarchical clustering selecting 31 clusters. The most common paths that people are taking are shown in Fig. 9.

The trajectory clusters give useful semantic information on the behaviour undertaken by the metro users. For instance, Cluster 14 indicates that most users enter the station by the north doors and go rather straight to the gates. Cluster 1 indicates that most people take also the north doors to exit the station. Cluster 25 represents fewer users exiting the gates and leaving the station to go through the south doors. Cluster 2 shows people with stationary activity near the gates. Cluster 26 indicates that few users buy a ticket before going through the gates. Cluster 5 indicates that after buying a ticket, users go straight to the gates to take the metro.



**Figure 6** Resulting partition of the CAVIAR data after running the relational analysis algorithm  
Properties of cluster 4 are given

The data processing of the Roma Underground (ATAC, Italy) corresponds to five hours and half of video, which is acquired at a rate of 5 frames/s. The same unsupervised clustering algorithm was applied to the data, which contains over 14 000 tracked objects. Fig. 10 presents the whole dataset of trajectories that we have analysed together with some of the clusters representative of the most common paths. For instance, trajectory cluster 20 in the upper right panel of the figure presents people entering the hall from the north doors. This is actually the biggest trajectory cluster followed by trajectory cluster 19 (lower left



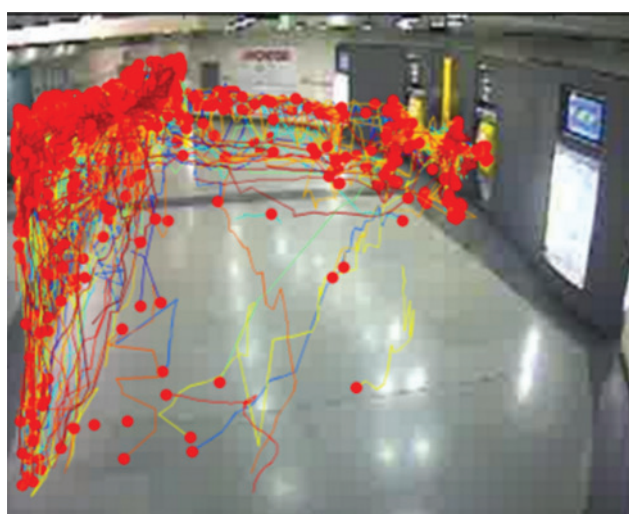
**Figure 7** Properties for cluster 5 in the CAVIAR data partition

**Table 4** Performance evaluation on the correspondence between the events annotated in the CAVIAR ground-truth and the relational analysis activity clusters

| Event               | Number of cases | Corresponding cluster | Detection performance |                   |                  |                    |
|---------------------|-----------------|-----------------------|-----------------------|-------------------|------------------|--------------------|
|                     |                 |                       | True positive, %      | False positive, % | True negative, % | Cluster quality, % |
| left baggage        | 4               | 4                     | 100                   | 0                 | 0                | 92                 |
| fight               | 12              | 0                     | 67                    | 0                 | 33               | 72                 |
| fight and fall down | 4               | 5                     | 75                    | 25                | 25               | 73                 |
| browsing            | 12              | 2                     | 86                    | 0                 | 14               | 74                 |

panel) showing people leaving the hall through north doors. Trajectory cluster 7 shows people exiting the hall through south doors; however, the number of elements of this cluster is smaller than the two previous ones. This means that south doors are less employed to enter/exit the hall.

We have further performed statistical analysis to measure the interactions between users and contextual objects. As mentioned in Section 3.2 there are two types of contextual objects of interest in the scene, the zones (mainly the hall), and the equipment (the gates and the vending machines). Over the whole observation time, concerning the Torino Underground, people were practically constantly present in the hall as we obtained a percentage of use of the hall of 91% of the observation time. The gates had a use of 36% indicating that the flow of people through the gates was not constant over the observation time. The two vending machines of the Torino station had a percentage of use of only 8% and 7%, respectively indicating that most people did not stop for a long-time or did not need to buy a ticket while in the station. This is further confirmed by the fact that most users buying tickets were detected as single

**Figure 8** Trajectories detected in one station of the Torino Underground (2025 trajectories during 50 min)

Beginnings of the trajectories are indicated by red points

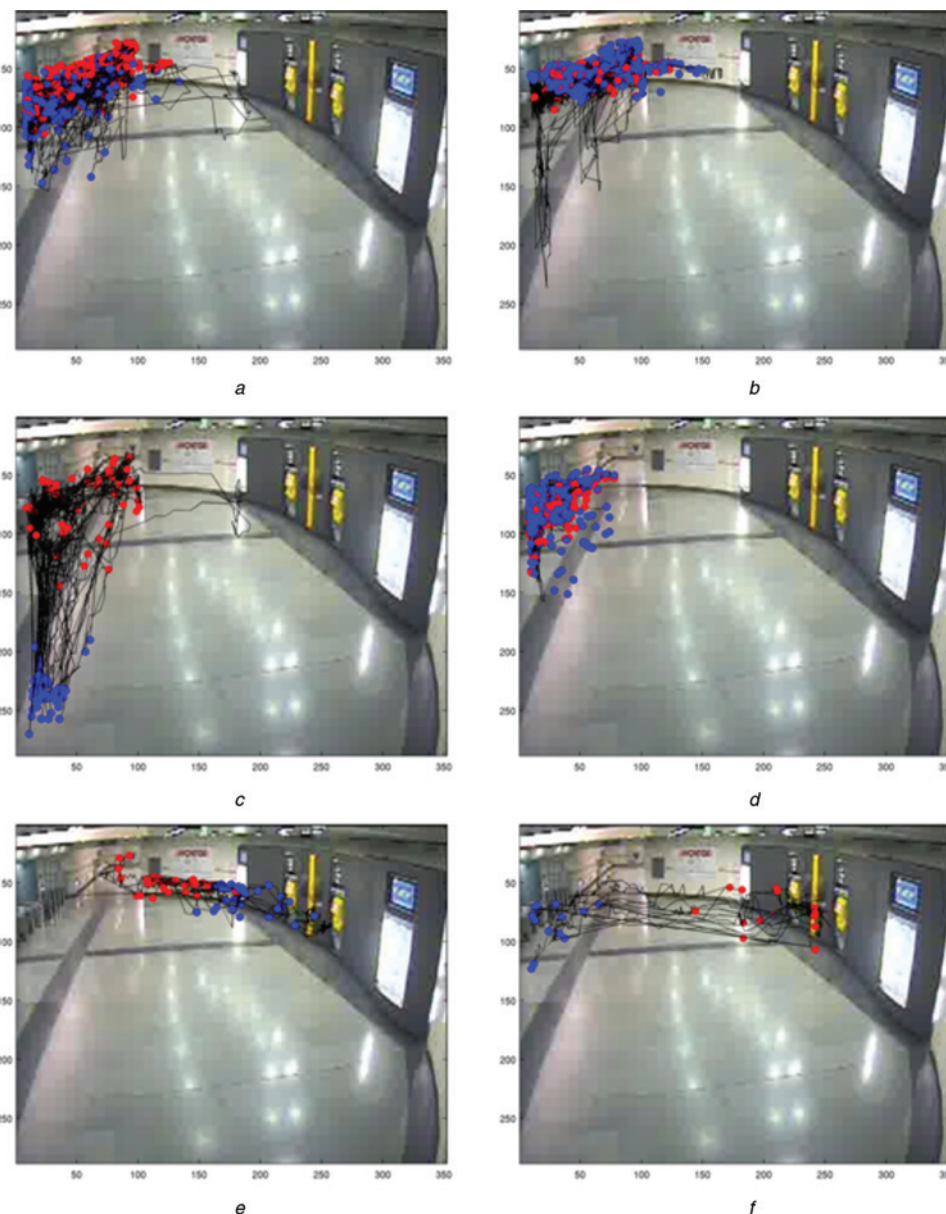
individuals and came more rarely to the machines as groups. No crowd (i.e. no queue) was detected at the vending machines or at the gates, whereas crowding was detected in the hall but at a low degree.

Regarding the interactions of people with contextual objects in the Roma station, people were observed in the hall 95% of the observation time and the gates were employed 91% of the time indicating that the flow of people through the gates is much more constant than, for instance, in the Torino Underground. The vending machines are often used 77% of the observation time. Compared with the Torino station, this utilisation time indicates that the Roma station is in general much busier. This is further confirmed by the fact that person groups are frequently detected at the vending machines (632 in the last 2 h), at the hall (1441) and at the gates (813). Crowding situations were detected only in the hall (227 crowd groups).

We have translated the information related to detected objects and statistical information in the format presented in section 4. As an example, a portion of the Torino semantic tables obtained is presented next (Tables 5–7).

From the contextual object table, we have been able to follow the evolution of the interactions with contextual objects. For Torino, Fig. 11 shows a graph on the evolution of the number of people present in the Torino hall during the observation time. For this observation period, the peak hour is detected at 6 h 45 min with 200 people in 5 min with an average of 180 people. Fig. 12 shows the evolution on the usage of a vending machine. Interestingly, a user spends more time (~40% increase of time) with the vending machine when the hall is not crowded.

For Roma, the number of people travelling through the hall is higher than in Torino with an average of 355 people constantly in the hall. In the last step towards knowledge discovery we have applied the relational analysis process explained before in Section 5. The input was the mobile object semantic table described above and containing in



**Figure 9** Clusters showing the most common paths obtained from the dataset shown in Fig. 7

*a* Cluster 14 with 443 trajectories

*b* Cluster 1 with 744 trajectories

*c* Cluster 25 with 50 trajectories

*d* Cluster 2 with 338 trajectories

*d* Cluster 26 with 41 trajectories

*e* Cluster 5 with 13 trajectories

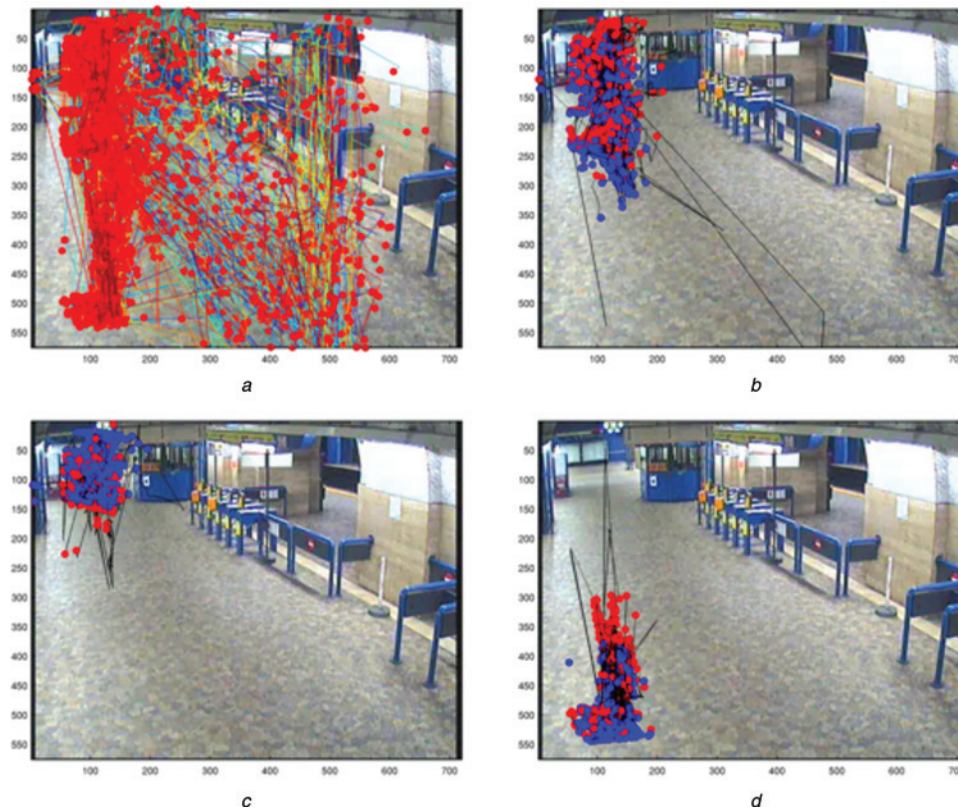
Beginnings of the trajectories are indicated by red points

Ends of trajectories are indicated by blue points

total 2025 detected objects in the case of the Torino data and 14757 for the Roma data. Some of the clusters returned are detailed next.

For the Torino dataset, the biggest cluster, labelled 0, contained 661 elements with the following variable properties: shape type: small person (100% of the elements); mobile object type: person (100% of the elements); duration: [0.04–5.64 s] (82% of the elements); significant event: inside\_zone\_hall (67% of the elements);

distance origin to destination: [0–59 cm] (47% of the elements); trajectory type: 1 (44% of the elements). In other words this activity cluster is made only of persons detected for a very short period of time (<6 s) and also moving in a short range of ~60 cm. They are associated with Trajectory type 1 'from gates to north doors' (shown in Fig. 9). Similar descriptions occur with the next biggest clusters but formed respectively of Unknown (524) and small PersonGroup (325) elements. Meaning that people take north doors as main exit; sometimes



**Figure 10** Clusters showing the most common paths obtained from the Roma dataset

*a* Whole dataset

*b* Cluster 20 with 1451 trajectories

*c* Cluster 19 with 809 trajectories

*d* Cluster 7 with 790 trajectories

Extremity points of the trajectories are represented as in Fig. 9 by red and blue points

groups of persons are formed in the flow of people but principally individuals are detected. This region (between gates and north doors) is the most visited of the scene.

The cluster labelled 4 is only made of persons (56 elements in the cluster) with the associated event 'stays at Gates'. They have been detected from  $\sim 6$ –15 s and by walking they are covering a distance between 60 cm and 2 m. The associated trajectory type is 14 (From north doors to gates). Thus, people entering the gates from north doors rarely have to wait at the gates to enter (relative small number of elements included in the cluster). This, however, may happen as 'gates–north doors' and vice versa is the main people path.

Another example of cluster activity is given by the cluster labelled 3, made only of Luggage elements (272 items) detected for a very short period of time (0.04–5 s) and mainly associated with trajectory type 2 'stationary at the gates' (shown in Fig. 9) meaning also that it is at the gates that most frequently people leave their luggage but only for a very short period of time. One last example of the Torino data is, for instance, Cluster 7, made up only of persons with associated significant event 'inside\_zone\_hall' and mainly trajectory type 25 'from gates to south doors' (Fig. 9). People following this path indeed walk longer to

go through the south doors and thus are longer inside the hall.

Because no ground-truth is available for the Torino data, it is not possible to perform the same evaluation of the clustering that we did for the CAVIAR dataset based on the calculation of true positives (TP) and false positives (FP). However, to have a quantitative measure that indicates the number of activities correctly detected by our clustering algorithm, we have considered TP detection, a cluster whose data correspond mainly after visual inspection by the end-user to a clear and meaningful description as for the clusters just presented above. FP is then defined as a cluster whose description cannot be associated with any coherent activity. In this sense, we had 20/27 clusters as TPs (activities corresponding to 90% of the detected objects) and 7/27 clusters as FPs.

For the Roma dataset, the biggest cluster contained 1672 elements. It is made only of small person groups detected inside the zone hall. They cover a distance of 50 cm to 1 m in 1–3 s. This means that although groups of persons appear very often from the north doors towards the hall (trajectory type 20; also shown in Fig. 10), they are able to walk at a normal speed. Another example for the Roma



**Table 5** Contextual objects semantic table

| Ctx_obj_type | Ctx_obj_type  | Startframe | Endframe          | Sig_event_type       | Rare_event_type             |
|--------------|---|------------|-------------------|----------------------|-----------------------------|
| 1            | platform  | 20050      | 92815             | inside_zone(14486)   | group_stays_inside_zone(30) |
| 2            | gates   | 20055      | 92745             | close_to(5103)       | group_stays_at(40)          |
| 3            | vendingmachine2   | 26560      | 90725             | close_to(1020)       | group_stays_at(200)         |
| 4            | vendingmachine1   | 30680      | 90650             | close_to(834)        | group_stays_at(160)         |
| Ctx_obj_id   | evnt_hist   |            |                   |                      |                             |
| 1            | insde_zone(14486) group_inside_zone(5523) crowd_inside_zone(1750) stays_inside_zone(1276) crowding_in_zone(109) |            |                   |                      |                             |
| 2            | close_to(5103) stays_at(3489) group_close_to(329) group_stays_at(40)  |            |                   |                      |                             |
| 3            | close_to(1020) stays_at(490) group_close_to(341) group_stays_at(200)  |            |                   |                      |                             |
| 4            | close_to(834) stays_at(482) group_close_to(307) group_stays_at(160)   |            |                   |                      |                             |
| Ctx_obj_id   | mob_obj_hist  |            | percent_of_use, % | mean_time_of_use, ms |                             |
| 1            | person(491) unknown(102) persongroup(136) luggage(34) crowd(4)  |            | 91.1296           | 35500.854            |                             |
| 2            | person(135) unknown(28) luggage(12) persongroup(17)   |            | 36.8284           | 35682.2396           |                             |
| 3            | unknown(19) person(15) luggage(10) persongroup(2)   |            | 8.8704            | 7401.087             |                             |
| 4            | unknown(12) person(11) luggage(8)   |            | 7.7504            | 16931.2903           |                             |

station is given by cluster labelled 6. This cluster (694 elements) is made up of only persons with associated significant event 'outside zone hall'. Indeed trajectory type 7 (Fig. 10) indicates that people are leaving the hall through the south doors. They cover a distance 1.67–2.30 m in ~3–4 s, which is a relative slow speed. The south doors

are not as busy as the north doors where numerous groups were detected also leaving the hall [cluster 5 (not shown) with 971 small groups detected].

For the Roma dataset, we have again considered TP detection as a cluster whose data correspond to a clear and

**Table 6** Events semantic table

| Evt_id | Evt_type     | Startframe | Endframe | Inv_objs_id | Ctx_type |
|--------|--------------|------------|----------|-------------|----------|
| 10161  | inside_zone  | 39085      | 39085    | 1573        | platform |
| 10162  | close_to     | 39085      | 39085    | 1444        | gates    |
| 10163  | inside_zone  | 39090      | 39090    | 1444        | platform |
| 10164  | inside_zone  | 39090      | 39090    | 1573        | platform |
| 10165  | close_to     | 39090      | 39090    | 1444        | gates    |
| 10166  | inside_zone  | 39095      | 39095    | 1444        | platform |
| 10167  | inside_zone  | 39095      | 39095    | 1573        | platform |
| 10168  | close_to     | 39095      | 39095    | 1444        | gates    |
| 10169  | inside_zone  | 39100      | 39100    | 1444        | platform |
| 10170  | inside_zone  | 39100      | 39100    | 1573        | platform |
| 10171  | close_to     | 39100      | 39100    | 1444        | gates    |
| 10072  | group_inside | 39105      | 39105    | 1573        | platform |

**Table 7** Mobile objects semantic table

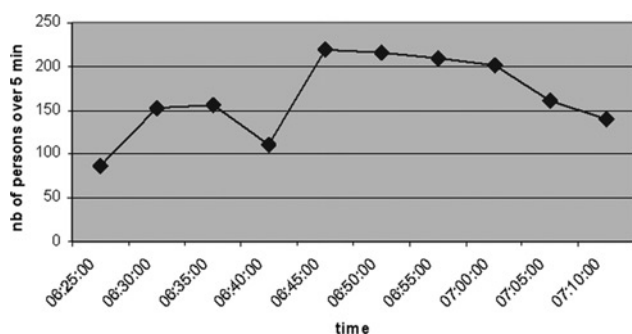
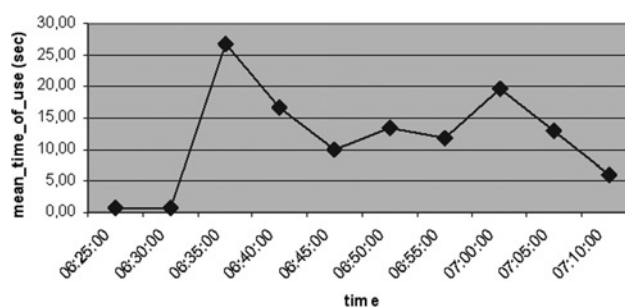
| Mob_obj_id | Mob_obj_type | Duration, s     | Traj_type | Dist_org_dest         | Shape_type         | Sig_event_id               |
|------------|--------------|-----------------|-----------|-----------------------|--------------------|----------------------------|
| 1          | person       | [48.24–52.04]   | 1         | [334.9776–485.3803]   | small person       | inside_zone_platform       |
| 2          | person       | [364.84–367.24] | 1         | [976.9466–1113.7527]  | small person       | stays_at_gates             |
| 3          | unknown      | [364.84–367.24] | 1         | [3158.1909–3230.0034] | small unknown      | void                       |
| 4          | luggage      | [295.24–298.84] | 2         | [671.8608–841.0493]   | small luggage      | void                       |
| 5          | personsgroup | [41.24–47.44]   | 14        | [60.075–205.4118]     | small personsgroup | group_inside_zone_platform |
| 6          | person       | [0.04–5.64]     | 1         | [334.9776–485.3803]   | small person       | inside_zone_platform       |
| 7          | unknown      | [33.84–40.24]   | 1         | [60.075–205.4118]     | small unknown      | void                       |
| 8          | person       | [0.04–5.64]     | 1         | [0–59.0762]           | small person       | inside_zone_platform       |
| 9          | personsgroup | [336.64–336.84] | 1         | [3630.5716–3723.4675] | small personsgroup | group_stays_at_gates       |

meaningful description, and FP as a cluster whose description cannot be associated with any coherent activity. In this sense, we had 72/109 clusters (activities) as TPs and 37/109 clusters as FPs.

Thus, this way the relational analysis can help us to group together people having similar behaviour. This is of particular interest to the end-users because the activities in the metro station can be better quantified.

In order to assess the impact of the object detection and tracking on the trajectory clustering step and how much the trajectory clustering process would then affect the results of the relational analysis, we evaluated our system on four different sets of tracked objects built from 3 h of observation of the Roma station. The first dataset (Dataset1: 2983 trajectories) contains all objects detected and tracked on this observation period including noisy data and fragmented trajectories. The next dataset (Dataset2: 2605 trajectories) contains all objects without trajectories of

short duration most likely to represent noise in the data. The third dataset (Dataset3: 1713 trajectories) is Dataset2 without all trajectories where the track for the last part of the trajectory was lost. Dataset4 (1713 trajectories) is Dataset 3 without all trajectories where the track was lost at some point and then continued. The lost part of the trajectory is then inferred. We evaluated the impact of the different trajectory datasets on the trajectory clustering tool by measuring the Silhouette, Dunn and Davis–Bouldin indexes. To evaluate how the relational analysis process is affected we calculated the resulting clustering quality as explained in Section 6. These results are shown in Table 8. As it can be observed, the quality of the trajectories extracted by the object detection and tracking module has an influence on the clusters obtained by the trajectory clustering module. The lesser the noise and lost tracks in the original dataset, the better the partition of the trajectory clusters. (The Dunn index monotonally increases with the trajectory quality. The Silhouette index generally increases and the Davis–Bouldin index generally decreases.) This influence is however less propagated into the relational

**Figure 11** Temporal evolution in the number of people occupying the station hall in the Torino station**Figure 12** Temporal evolution of the mean time a user spends on a vending machine in the Torino station hall

**Table 8** Evaluation of the impact of the object detection and tracking on the trajectory clustering and relational analysis processes

| Input    | Number of trajectories | Trajectory performance indexes |       |               | Relational analysis performance indexes |                              |
|----------|------------------------|--------------------------------|-------|---------------|---|------------------------------|
|          |                        | Silhouette                     | Dunn  | Davis–Bouldin | Clustering quality, %                   | Resulting number of clusters |
| dataset1 | 2983                   | 0.314                          | 0.049 | 15.571        | 70                                      | 35                           |
| dataset2 | 2605                   | 0.254                          | 0.056 | 13.79         | 69                                      | 46                           |
| dataset3 | 1713                   | 0.325                          | 0.071 | 17.632        | 69                                      | 42                           |
| dataset4 | 1476                   | 0.343                          | 0.091 | 17.854        | 69                                      | 39                           |

analysis as the quality of the final activity clusters remains rather constant.

Regarding our implementation, 1 h of video takes  $\sim 1$  min to be processed off-line by the trajectory clustering, statistical and relational analysis. This is a reasonable processing time in adequacy with end-user requirements.

## 8 Discussion and conclusion

In this paper we have presented how knowledge discovery can be achieved on large recordings of video using an efficient knowledge representation format. The richness of the representation comes from the fact that both moving objects and the contextual objects from the scene are studied together with their interaction. Yet, the proposed representation provides a useful support and enables all activity knowledge to be structured into three different appropriate tables, namely mobile objects, contextual objects and video events. The proposed representation supports a rich set of spatial topological and temporal relations and captures not only quantitative properties but also higher semantic concepts. Furthermore, a first layer of meaningful knowledge is directly extracted from the video streams by detecting events corresponding to the interactions between moving objects and contextual objects. A second layer of knowledge is extracted by the off-line long-term analysis of these interactions. First, statistical information is obtained from the mobile objects (in particular, their trajectory) and the contextual objects as well as their interactions (i.e. events). Statistical information is a major information source for the end-user. For instance, on large metro video recordings the average number of people localised in specific zones of interest or interacting with particular equipment and its evolution with time provide operational information on how to manage the metro station on a day-by-day basis. We are currently analysing sequentially chunks of video of  $\sim 2$  h and then will further analyse the temporal evolution for durations such as 1 day or 1 week. On a second step, we perform the trajectory characterisation by employing a hierarchical algorithm. It must be remarked that there is a scalability

issue when employing standard hierarchical algorithms, and these may not work efficiently for very large datasets. An alternative solution is to implement, for instance, parallel hierarchical clustering algorithms [45]. However, to solve the main problem concerning accessing and processing a larger number of stored data, we are currently working on a new version where we will be able to update on-line the clusters for continuous processing. In this new version also other features such as duration, mean speed and distance walked are currently being studied in the clustering of trajectories.

The semantic knowledge gained from trajectory characterisation and statistical analysis is then used for the discovery of complex relationships. The relational analysis proposed in this paper shows how to highlight hidden relations between people, their trajectories (behavioural information) and the significant interactions between themselves and contextual objects. Thus, relational analysis can define the typical activities in the subway represented as activity clusters.

The performance evaluation has been performed at the knowledge discovery level by providing manually ground-truth trajectories in two sites. In a building hall for the CAVIAR project and in the Torino Metro hall for the CARETAKER project. Performance evaluation is a sensitive point in the knowledge discovery because of the large number of parameters to employ and the small description on expected results to be provided by the end-users. In opposition to what is usually done in trajectory and activity clustering, we have proposed a new framework for evaluation. With this framework we have been able to assess the system. These results are however to be taken with care as the size of the data analysed is small. This unfortunately has been an endemic problem in the area of computer vision, where annotated datasets (with an evaluation ground-truth) are difficult to find. More work is still to be done, not only for evaluating more hours of video, but also, for taking into account other types of parameters such as spatial and temporal granularity of the scene.

## 9 References

- [1] MOESLUND T., HILTON A., KRÜGER V.: 'A survey of advances in vision-based human motion capture and analysis', *Comput. Vis Image Underst.*, 2006, **104**, (2–3), pp. 90–126
- [2] ZHANG H., KANTANKANHALLI A., SMOLIAR S.: 'Automatic partitioning of full-motion video', *ACM Multimed. Syst.*, 1993, **1**, (1), pp. 10–28
- [3] TONG S., CHANG E.: 'Support vector machine active learning for image retrieval'. Proc. ACM Multimedia, Ottawa, Canada, 2001, pp. 107–118
- [4] CHONG-WAH N., TING-CHUEN P., HONG-JIANG Z.: 'Video processing: on clustering and retrieval of video shots'. Proc. 9th ACM Int. Conf. Multimedia MULTIMEDIA'01, 2001
- [5] EWERTH R., FREISLEBEN B.: 'Semi-supervised learning for semantic video retrieval'. Proc. 6th ACM Int. Conf. Image and Video Retrieval CIVR '07, 2007
- [6] SMEATON A.F.: 'Techniques used and open challenges to the analysis, indexing and retrieval of digital video', *Inf. Syst.*, 2007, **32**, (4), pp. 545–559
- [7] LE T.-L., BOUCHER A., THONNAT M.: 'Subtrajectory-based video indexing and retrieval'. Proc. 13th Int. Conf. Advances of Multimedia Modelling, 2007, vol. 1, pp. 418–427
- [8] VELASTIN S.A., BOGHOSSIAN B.A., LAI LO B.P., SUN J., VICENCIO-SILVA M.A.: 'PRISMATICA: toward ambient intelligence in public transport environments', *IEEE Trans. Syst. Man Cybern. A*, 2005, **35**, pp. 164–182
- [9] PIATER J., RICHTETTO S., CROWLEY J.: 'Event based activity analysis in live video using a generic object tracker'. Proc. 3rd IEEE Workshop on Performance Evaluation of Tracking and Surveillance (PETS), Copenhagen, Denmark, 2002
- [10] CUPILLARD F., BRÉMOND F., THONNAT M.: 'Automatic visual recognition for metro surveillance'. Proc. Int. Conf. Measuring Behavior, Wageningen. The Netherlands, 2005
- [11] MARCOTORCHINO F., MICHAUD P.: 'Optimisation en analyse ordinaire des données', Masson 1978
- [12] <http://homepages.inf.ed.ac.uk/rbf/CAVIAR/>
- [13] PATINO J.L., BENHADDA H., CORVEE E., BREMOND F., THONNAT M.: 'Video-data modelling and discovery'. IET Int. Conf. Visual Information Engineering VIE'07, 2007, p. 9
- [14] GEORIS B., BREMOND F., THONNAT M., MACQ B.: 'Use of an evaluation and diagnosis method to improve tracking performances'. Proc. 3rd IASTED Int. Conf. Visualization, Imaging and Image Process. (VIIP), 2003, p. 2
- [15] WANG L., HU W., TAN T.: 'Recent developments in human motion analysis', *Pattern Recognit.*, 2003, **36**, pp. 585–601
- [16] FUSIER F., VALENTIN V., BRÉMOND F., ET AL.: 'Video understanding for complex activity recognition', *Mach. Vis. Appl. J.*, 2007, **18**, pp. 167–188
- [17] CUPILLARD F., BRÉMOND F., THONNAT M.: 'Tracking group of people for video surveillance'. IEEE Proc. 2nd European Workshop on Advanced Video-Based Surveillance System AVBS'02, 2002
- [18] AVANZI A., BREMOND F., TORNIERI C., THONNAT M.: 'Design and assessment of an intelligent activity monitoring platform'. EURASIP, 2005, pp. 2359–2374
- [19] AVANZI A., BRÉMOND F., THONNAT M.: 'Tracking multiple individuals for video communication'. Proc. IEEE Int. Conf. Image Processing, 2001
- [20] VU V.T., BREMOND F., THONNAT M.: 'Automatic video interpretation: a novel algorithm for temporal scenario recognition'. Proc. 18th Int. Joint Conf. Artificial Intelligence. IJCAI'03, 2003, pp. 1295–1302
- [21] <http://www.cvg.rdg.ac.uk/PETS2006/index.html>
- [22] PICIARELLI C., FORESTI G.L., SNIDARO L.: 'Trajectory clustering and its applications for video surveillance'. IEEE Int. Conf. Advanced Video and Signal Based Surveillance, AVSS'05, 2005
- [23] ANJUM N., CAVALLARO A.: 'Single camera calibration for trajectory-based behavior analysis'. IEEE Int. Conf. Advanced Video and Signal Based Surveillance, AVSS'07, 2007
- [24] NAFTEL A., KHALID S.: 'Classifying spatiotemporal object trajectories using unsupervised learning in the coefficient feature space', *IEEE Trans. Multimed. Syst.*, 2006, **12**, (3), pp. 45–52
- [25] ANTONINI G., THIRAN J.P.: 'Counting pedestrians in video sequences using trajectory clustering', *IEEE Trans. Circuits Syst. Video Technol.*, 2006, **16**, (8), pp. 1008–1020
- [26] CALDERARA S., CUCCHIARA R., PRATI A.: 'Detection of abnormal behaviors using a mixture of Von Mises distributions'. IEEE Int. Conf. Advanced Video and Signal Based Surveillance, AVSS'07, 2007
- [27] GAFFNEY S., SMYTH P.: 'Trajectory clustering with mixtures of regression models'. Proc. Int. Conf. Knowledge Discovery and Data Mining, CA, USA, 1999

- [28] PORIKLI F.: 'Learning object trajectory patterns by spectral clustering'. Proc. IEEE Int. Conf. Multimedia and Expo ICME '04, 2004, vol. 2, pp. 1171–1174
- [29] BASHIR F.I., KHOKHAR A.A., SCHONFELD D.: 'Object trajectory-based activity classification and recognition using hidden Markov models', *IEEE Trans. Image Process.*, 2007, **16**, (7), pp. 1912–1919
- [30] OLIVER N.M., ROSARIO B., PENTLAND A.P.: 'A Bayesian computer vision system for modeling human interactions', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2000, **22**, (8), pp. 831–843
- [31] JIAN F., WU Y., KATSAGGELOS A.: 'Abnormal event detection from surveillance video by dynamic hierarchical clustering'. ICIP, 2007
- [32] XIANG T., GONG S.: 'Video behaviour profiling and abnormality detection without manual labelling'. Proc. IEEE Int. Conf. Computer Vision, 2005, pp. 1238–1245
- [33] XIANG T., GONG S.: 'Incremental and adaptive abnormal behaviour detection', *Comput. Vis. Image Underst.*, 2008, in press
- [34] TOSHEV A., BRÉMOND F., THONNAT M.: 'An a priori-based method for frequent composite event discovery in videos'. Proc. 2006 IEEE Int. Conf. Computer Vision Systems, 2006
- [35] HARTIGAN J.A.: 'Clustering algorithms' (John Wiley & Sons, New York, 1975)
- [36] KOHONEN T.: 'Self-organizing maps' (Springer-Verlag, New York, 2001)
- [37] KAUFMAN L., ROUSSEEUW P.J.: 'Finding groups in data. An introduction to cluster analysis' (Wiley, New York, 1990)
- [38] CAMPELLO R.J.G.B., HRUSCHKA E.R.: 'A fuzzy extension of the silhouette width criterion for cluster analysis', *Fuzzy Sets Syst.*, 2006, **157**, pp. 2858–2875
- [39] DUNN J.: 'Well separated clusters and optimal fuzzy partitions', *J. Cybern.*, 1974, **4**, pp. 95–104
- [40] DAVIES D.L., BOULDIN D.W.: 'A cluster separation measure', *IEEE Trans. Pattern Recognit. Mach. Intell.*, 1979, **1**, pp. 224–227
- [41] BENHADDA H., MARCOTORCHINO F.: 'L'analyse relationnelle pour la fouille de grandes bases de données', *Revue des Nouvelles Technologies de l'Information*, 2007, **RNTI-A-2**, pp. 149–167
- [42] BENHADDA H.: 'La similarité régularisée et ses applications en classification automatique, PhD thesis, Paris VI University, 1998
- [43] BENHADDA H., MARCOTORCHINO F.: 'Introduction à la similarité régularisée en analyse relationnelle', *Revue de Statistique Appliquée*, 1998, **46**, (1), pp. 45–69
- [44] MARCOTORCHINO F.: 'Seriation problem: an overview', *Appl. Stochastic Models Data Anal.*, 1991, **7**, pp. 139–151
- [45] RAJASEKARAN S.: 'Efficient parallel hierarchical clustering algorithms', *IEEE Trans. Parallel Distrib. Syst.*, 2005, **16**, pp. 497–502