



Robust abandoned object detection integrating wide area visual surveillance and social context [☆]

James Ferryman ^{a,*}, David Hogg ^b, Jan Sochman ^c, Ardhendu Behera ^b, José A. Rodriguez-Serrano ^d, Simon Worgan ^e, Longzhen Li ^a, Valerie Leung ^g, Murray Evans ^a, Philippe Cornic ^f, Stéphane Herbin ^f, Stefan Schlenger ^h, Michael Dose ^h

^a Computational Vision Group, School of Systems Engineering, University of Reading, RG6 6AY, UK

^b School of Computing, University of Leeds, LS2 9JT UK

^c Centre for Machine Perception, Department of Cybernetics, Faculty of Electrical Engineering, Czech Technical University, Karlovo namesti 13, 121 35 Praha 2, Czech Republic

^d Xerox Research Centre Europe, 6 Chemin de Maupertuis, 38240 Meylan, France

^e formerly University of Leeds

^f Department of Information Processing and Modelling, ONERA, BP 80100, 91123 Palaiseau Cedex, France

^g MathWorks, Les Montalets, 2 rue de Paris, 92190 Meudon, France

^h L-1 Identity Solutions, Universitaetsstr.160, 44801 Bochum, Germany

ARTICLE INFO

Article history:

Available online 4 February 2013

Keywords:

Wide area video surveillance
Behaviour analysis
Abandoned objects

ABSTRACT

This paper presents a video surveillance framework that robustly and efficiently detects abandoned objects in surveillance scenes. The framework is based on a novel threat assessment algorithm which combines the concept of ownership with automatic understanding of social relations in order to infer abandonment of objects. Implementation is achieved through development of a logic-based inference engine based on Prolog. Threat detection performance is conducted by testing against a range of datasets describing realistic situations and demonstrates a reduction in the number of false alarms generated. The proposed system represents the approach employed in the EU SUBITO project (Surveillance of Unattended Baggage and the Identification and Tracking of the Owner).

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

In recent years there have been a number of incidents where terror organisations have planted explosive devices in ordinary baggage to cause immense disruption in mass transportation networks and other areas of critical infrastructure. Due to the potentially devastating consequences of such terrorist activity, the monitoring and surveillance of unattended baggage has become a priority for the security operators of mass transportation networks and other critical infrastructure. The overriding goal is to minimise the number of false alarms. Towards this goal, the main contribution of this work is the development and evaluation of behaviour analysis methodology permitting robust identification of a baggage-owner while minimising false positives. The approach taken advances the state of the art in abandoned bag detection by introducing the concept of ownership and combines it with automatic understanding of social groups to infer abandonment. To achieve the goal, a framework (see Fig. 1) has been developed consisting of a complete fourfold process, detection – tracking – situation

analysis – threat assessment. This paper is divided as follows. Firstly, in Section 2 related research is detailed, followed in Sections 3–5 by descriptions of the system components. In Section 6 the datasets used and results of experiments are presented before concluding in Section 7 with conclusions and recommendations for future research.

2. Related work

There exists a significant body of academic research addressing the task of robustly identifying abandoned baggage in public spaces. Most authors treat detection of abandoned (or left) objects, especially luggage, as the task of static object detection, with (Birch et al., 2011; Tian et al., 2010) or without (e.g. (Evangelio and Sikora, 2011; Porikli et al., 2008)) the application of tracking. Tian et al. (2010) present a framework to detect abandoned and removed scene objects based on background subtraction and foreground analysis, combined with tracking output to reduce false positives. Birch et al. (2011) employ motion segmentation based on a GMM with fast learning and a motion history image (MHI). For tracking of stationary objects, the edge map (3×3 Sobel filter) for each pixel is computed and matched) by correlation of edge directions. A comparative evaluation of stationary foreground detection

[☆] This document is a collaborative effort.

* Corresponding author. Tel.: +44 118 3786697; fax: +44 118 9751994.

E-mail addresses: james@computer.org, j.m.ferryman@reading.ac.uk (J. Ferryman).

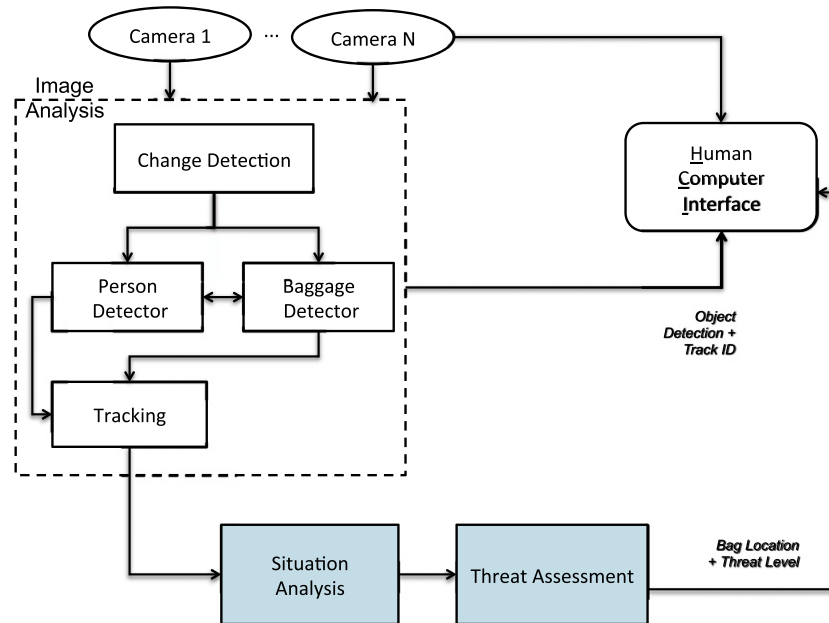


Fig. 1. General framework of the automated threat detection system.

algorithms based on background subtraction is given in (Bayona et al., 2009).

There has been some attempt at human activity recognition and association to scene objects. In (Lu et al., 2007) moving objects are tracked using shape and colour features and Kalman-based filtering, and classified using eigen features and Support Vector Machine. A package is defined as a non-human object and package ownership analysis performed using HMM-based human activity recognition.

2.1. Dataset based challenges

The most widely used datasets with which to evaluate approaches to abandoned bag detection have been from (PETS, 2007; PETS, 2006) and from the UK Home Office i-LIDS (2007). The dataset provided for the PETS (2006) challenge consists of 7 multi-camera scenarios involving an increasing number of people and passers-by. Most of the submissions to PETS2006 were based on background subtraction combined with a blob tracker (Auvinet et al., 2006; Guler and Farrow, 2006; Krahnstoever et al., 2006; Li et al., 2006; Martínez-del-Rincón et al., 2006; Smith et al., 2006), with the exception of Lv et al. (2006) who rely on a more realistic human model by incorporating a human detector. Most often, when an object is not moving and its size is beneath a given threshold, it is assumed to be a standing bag. Smith et al. (2006) propose a probabilistic approach in which people and bags are classified based on the immediate history of their size and velocity. Another approach from PETS2006 is to use a slow-decay background model to detect stationary objects (Guler and Farrow, 2006). To be able to apply the PETS2006 rules for abandoned baggage (the owner is further than a metres for more than b seconds), the owner is usually defined as the nearest tracked object when the standing bag appears (Krahnstoever et al., 2006; Lv et al., 2006) or by examining blob splits during tracking (Auvinet et al., 2006; Guler and Farrow, 2006; Smith et al., 2006). When a standing bag and its owner are identified, it is straightforward to apply the PETS2006 abandoned-bag rules. The simplicity of the scenarios allows very limited situation awareness and was designed mainly to test if the low level processing stages are sufficient to cope with real-world scenarios.

The PETS (2007) challenge focusses on two additional scenarios: theft and loitering. The videos are much more challenging from the

tracking point of view as the scenes are more crowded. There are eight scenarios, each viewed from four cameras. Two submissions to the challenge go beyond classical approaches to blob tracking and split-track analysis (such as (Arsić et al., 2007; Dalley et al., 2007)) and slowly/quickly adapting background models (such as Porikli and Yin (2006)). Firstly, Ribeiro et al. (2007) use a Temporal-JointBoost algorithm for each blob being tracked to classify it into a person-walking, not moving, a person picking-up/leaving a bag, or an abandoned bag. The basic idea is to incorporate temporal features (optical flow, motion energy) into the classification process over some temporal window. Secondly, Ardö and Aström (2007) use an HMM to improve the temporal consistency of the tracking and show how to use an HMM efficiently in this setting. These approaches demonstrate the potential advantages of considering a longer temporal window for activity analysis. Nevertheless, the situation awareness in the PETS 2007 challenge is again very simple – reduced to comparing the distance of a bag to its owner (abandoned bag, theft) or measuring the time for which a person stays in the scene (loitering).

The UK Home Office have developed an image library (i-LIDS, 2007) to help researchers and designers to evaluate video based detection systems to meet Government requirements. The i-LIDS library includes an abandoned luggage dataset including several challenges of single instances of left luggage on a metro platform in the presence of passing passengers and trains. While the dataset is useful for evaluating detection algorithms it remains limited because it is monocular and also does not contain examples of specific behavioural interactions.

2.2. Limitations of existing approaches

It is clear that a global analysis of the situation rather than just examining each agent's behaviour independently, would be beneficial in many situations. The motivation for this is illustrated by a scenario similar to that of (PETS, 2007) where a family or a group of friends comes together and one of them leaves his/her bag with the others. Any threat detection system treating the individuals independently would inevitably report an abandoned bag, as the criteria specified in (PETS, 2006) that the bag is abandoned if the owner is further than a metres for more than b seconds, is fulfilled. For treating these more complex scenarios, the approaches described

above may be insufficient and it may be necessary to derive a more complete activity analysis. A significant corpus of the computer vision and artificial intelligence literature attacks the problem of understanding activities from visual input. While logic and grammar-based representations, with or without combination with statistical approaches, (Hongeng et al., 2004; Ivanov and Bobick, 2000; Joo and Chellappa, 2006; Shet et al., 2005) organise knowledge in a flexible, powerful and clean way, one drawback of these approaches is that they are unable to propagate the uncertainty in the primitive detections. Hidden Markov Models (Brand et al., 1997) and other flavours of dynamic Bayesian network provide a powerful generalisation of stochastic finite state automata to deal with such uncertainty. Another related approach is the so-called propagation network (Shi et al., 2004). In recent work, Damen and Hogg (Damen and Hogg, 2012) first specify activities using a multi-set attribute grammar and then convert it to an equivalent Bayesian network. A more general tool which converts first-order logic predicates into an equivalent Bayesian network is the framework of Markov logic networks (Richardson and Domingos, 2006), which have also been applied to activity analysis (Tran and Davis, 2008). An entirely different approach is to detect events from image pixels directly rather than by reasoning about the interactions between specific agents, for instance (Li et al., 2008; Wang et al., 2009).

Whilst these approaches are easily configured to output whether an activity is normal or abnormal, they lack the explanatory power of grammar and logic-based methods (i.e. why it is abnormal).

None of the approaches described in the literature, however, have combined the concept of ownership with recognition of social groups, to reduce the number of false positives in detection of abandoned objects.

3. Object detection and tracking

The framework, shown in Fig. 1, supports application of a range of object detectors and trackers including the POM person detection method of Berclaz et al. (2009) and tracking-by-detection of Breitenstein et al. (2011), both of which operate at low frame rates (2–4 fps) or offline. While detection and tracking is not the main contribution of this paper, brief descriptions are given to methods which have been developed to permit the overall framework to operate online and with multiple cameras.

3.1. Baggage detection

Baggage hypothesis generation is based on static change detection using the dual background approach of Porikli et al. (2008)

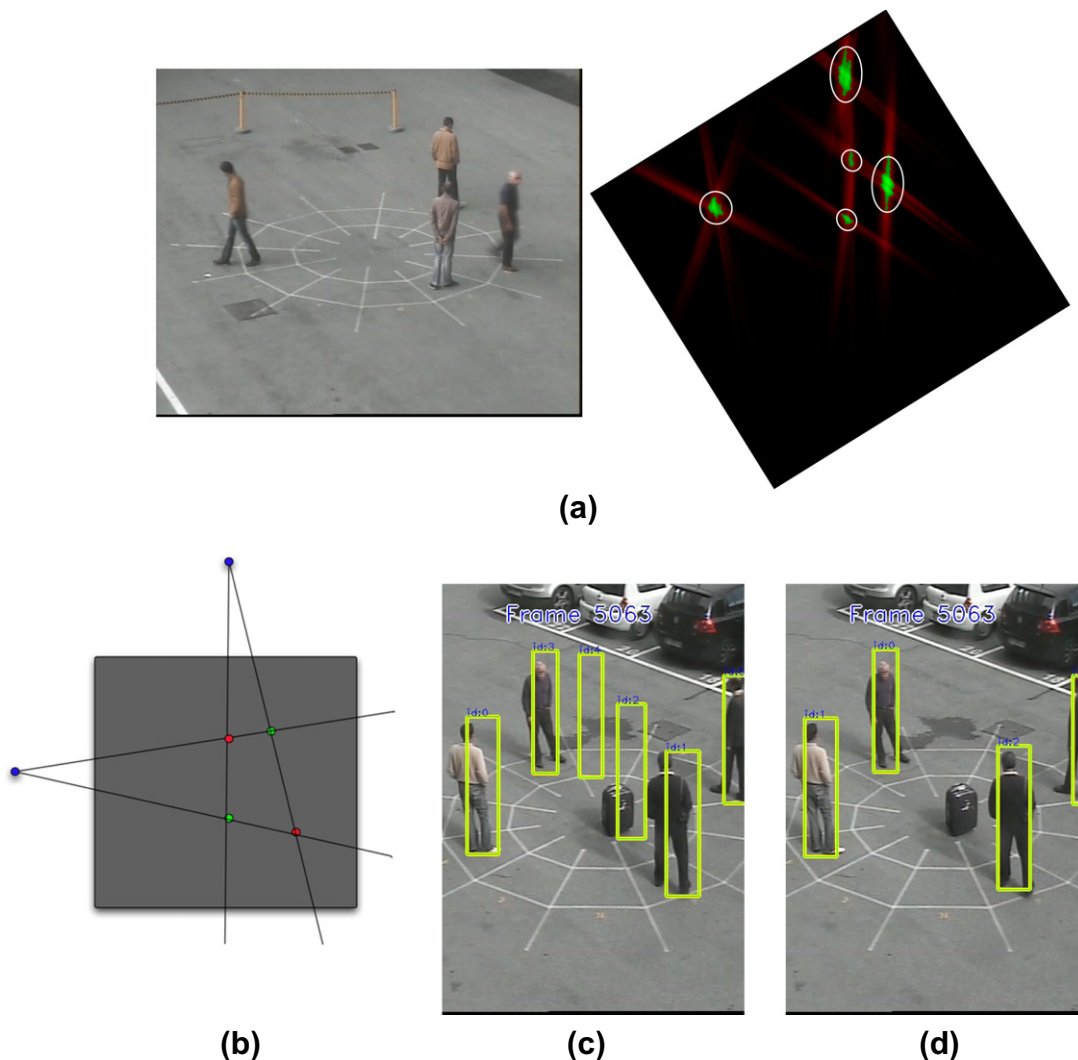


Fig. 2. Synergy map: (a) Detection of all pedestrians requires a threshold on synergy map to be set to value that permits ghost detection to pass through. (b) Ghost positions (red) can be predicted if correct positions (green) are known or can be estimated. (c–d) Bounding boxes resulting from detections without (c) and with (d) ghost prediction and suppression, for the same frame of video. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article).

adapted to use the efficient implementation of the Gaussian Mixture Model in (Zivkovic, 2004). Bag verification consists of application of a combination of filters including both 2D and 3D geometric filters, foreground/background similarity filter, and temporal filtering to check for persistence of the static regions.

3.2. Person detection

Person detection is based on the homography based multi-camera approach of Yildiz and Akgul (2010), extended with a novel approach for ghost suppression. First, a synergy map, the result of projecting detected foreground from each camera view to a single plane, is created, as shown in Fig. 2. In practise, the reverse process is used with sampled cells on the synergy map, each corresponding to a vertical cuboid in space of fixed person height, back-projected to the bounded rectangles in the original images. The process is applied for an image resolution-limited "infinite" number of planes in

a very efficient and fully real-time manner without hardware acceleration.

For a given location (x,y) in the synergy map (which corresponds to a small rectangular region on the ground plane), the value $S(x,y)$ accumulating the evidence of a person's presence can be calculated as:

$$S(x,y) = \frac{1}{|I|} \sum_{i \in I} \frac{\sum_{u=u_0}^{u_1} \sum_{v=v_0}^{v_1} p(u,v,i)}{A(Z(x,y,i))} \quad (1)$$

where I is the set of images into which the cuboid can be visibly projected, $Z(x,y,i) = \{(u_0, v_0), (u_1, v_1)\}$ is the bounding box projection of the cuboid corresponding to a specific synergy map pixel (x,y) into image i as defined by two extreme corner points. $A(s)$ is a function to calculate the area of any shape s , and

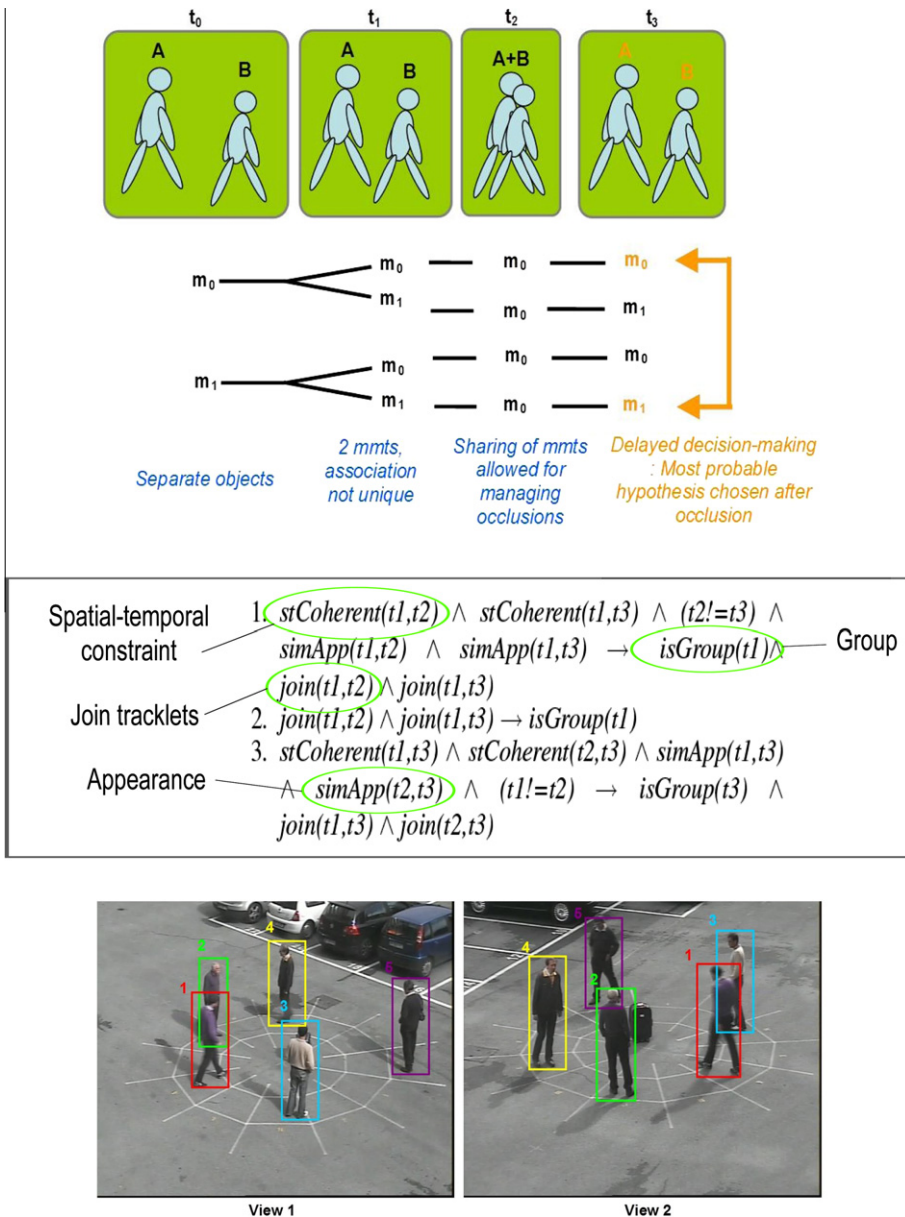


Fig. 3. Tracking processes. Top: illustrating how measurement-sharing in video-MHT overcomes short-term occlusions. Middle: examples of tracklet association rules used in the MLN formalism. Spatial-temporal coherence and appearance information are used as inputs. The inference of groups and the joining of tracklets are two of the outputs. Bottom: example tracking output for two cameras showing objects IDs.

$$p(u, v) = \begin{cases} 1, & \text{if } I(u, v) \text{ is foreground} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

Candidate objects are represented by peaks in the synergy map, obtained via thresholding. Ghost detections can occur where lines from different cameras to different objects intersect. To prevent ghosts becoming new tracking targets, a suppression map is generated in the regions of high ghost probability and subtracted from the synergy map. Frame-to-frame tracking of peaks further reinforces probable objects' location.

3.3. Tracking

A multi hypothesis tracker is used (Blackman, 2004) modified for application to tracking of extended objects. First, to handle short-term occlusions and the merging of measurements from different persons in the detection process, measurement-sharing between track hypotheses is allowed. This concept is illustrated in Fig. 3 (Top). Secondly, the measurement-to-track association cost is modified to allow image features, specifically two hue-saturation histograms corresponding to the top and bottom halves of a person, to be used in addition to a simple Brownian motion model. Each model is updated using the exponentially weighted moving average (EWMA). The association score between a predicted state and a measurement is a product of the normalised histogram intersection distance between their histograms and the normalised Euclidean distance between their positions in 3D.

To overcome track fragmentations caused by long-term or complex patterns of interaction between people, long term tracking based on tracklet association is used. The approach is based on a Markov Logic Network (MLN) (Leung and Herbin, 2011) where the notion of a group to account for generic interaction between people is introduced. The scores for possible associations are calculated using both spatial-temporal constraints and appearance information. Associations are not only considered for tracklets that can be directly joined together; but are extended to tracklets separated by a group in space and time. It therefore handles the formation and splitting of groups, reducing track fragmentations and allowing longer tracks to be formed. Examples of the tracklet association rules are shown in Fig. 3 (Middle) and example final tracking output in Fig. 3 (Bottom).

4. Situation analysis

Situation analysis is an intermediate step towards threat assessment and is defined as the description of the relationships between people and bags that can be inferred from the behaviour of the participating agents. This contribution focuses on two kinds of relationship: who owns each bag, and who knows who. The analysis takes object tracks and class information as input and describes the state of the world (i.e. the scene) in terms of the observed agents and their behaviour. The following stage (threat assessment) determines whether the state of the world constitutes a possible threat (i.e. there is a truly abandoned bag). The main contribution is the combination of the automatic understanding of social relationships with the concept of ownership to reduce the number of false alarms.

4.1. Bag ownership

For the reported experiments in this paper, a bag is detected when it appears stationary in the scene, having been placed there by a person. At this stage, detection of a bag as it is carried into or out of the scene has not been incorporated. The ownership of each bag is inferred by simply looking for a person in the proximity of the bag over a fixed time interval prior to its appearance. The person is also required to be stationary at the time the bag-drop is

hypothesised to occur. Specifically, in the experiments reported here, any person is assumed to be an owner if they are temporarily stationary within one metre of the bag at any point within one second prior to its appearance. Note that multiple possible owners are allowed, not because this is expected to be the case in reality but in order to reduce false alarms through taking both hypotheses through into the threat assessment.

4.2. Inference of social relations

Social groups are a very common phenomena in human crowds, with empirical studies suggesting that about 74% of people come in a group to a social event (Aveni, 1977) and about 50–70% (depending on the environment) are in a group during casual walking (Rudloff et al., 2011). Despite this high percentage, the prevailing crowd behaviour models in today's simulation tools (Challenger et al., 2009), computer graphics applications (Reynolds, 1987) and in particular in activity recognition and computer vision (PETS, 2006) are based on modelling each individual independently. An online algorithm has been developed for automatic detection of social groups within crowds, based on the analysis of the way the social relations influence the walking behaviour of the group members.

The method is based on the social force model (SFM) (Helbing and Molnar, 1995; Moussaïd et al., 2010) widely used in the crowd simulation community. In this, each individual's movement is influenced by notional forces operating between individuals. Depending on whether two individuals (a) know each other or (b) do not know each other, the Social Force Model produces different sets of trajectories for these individuals. Until recently, these attempts were based on human designed forces without proper evaluation. Only recently, the model has been calibrated on real-world video sequences resulting in a model that realistically predicts avoidance behaviour of a walking group (Moussaïd et al., 2009; Singh et al., 2009) and later in a model with all its parameters, including group behaviour, estimated from real data (Moussaïd et al., 2010).

The method employed in this work solves the inverse problem: knowing the trajectories, what are the social forces, and thus the relations, that caused that behaviour. The method is used in the framework to infer the social relations between the individuals in a scene and thereby to inform threat assessment as explained in Section 5.

The authors are aware of only two approaches aiming explicitly at social group inference (Ge et al., 2009; Jacques et al., 2007) and one paper using social groups to improve tracking (Pellegriani et al., 2010). In (Jacques et al., 2007) the groups are detected when two individuals keep close enough for a significant fraction of time over a given period. Experiments undertaken by the authors have shown that such simple measures are not sufficient for reliable group inference in complex scenes. In the proposed approach the calibrated SFM instead is relied upon. Similar measurements were used in (Pellegriani et al., 2010) to improve tracking by jointly tracking and inferring the social groups.

Also based on distance, but including the difference in velocity as well as position, the method proposed in (Ge et al., 2009) applies clustering to the (complete) person trajectories. The merging criterion takes into account the fraction of time in which the individuals are seen close to each other and allows the addition of a person to the group only if they have been close to at least half of its members. Fig. 4 illustrates the Social Force Model. Full details of the approach are given in (Sochman and Hogg, 2011).

5. Threat assessment

The threat assessment stage determines whether the inferred situation constitutes a threat, utilising the inferred knowledge of

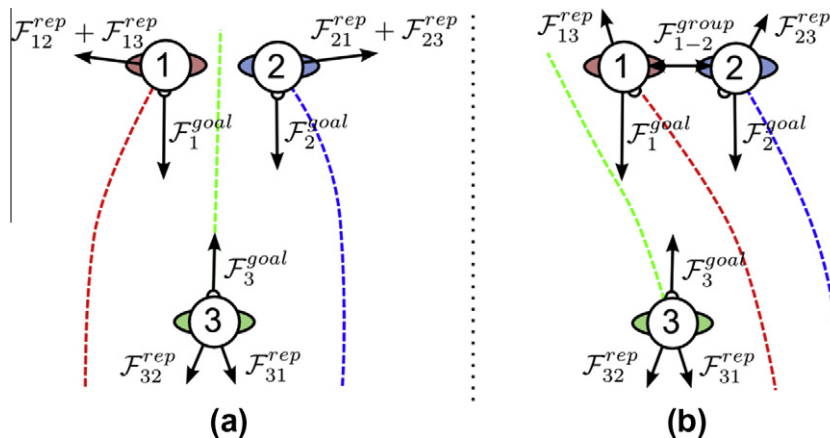


Fig. 4. Depending on whether the individuals 1 and 2 (a) do not know each other or (b) know each other, the Social Force Model produces different sets of trajectories combining together repulsive (\mathcal{F}^{rep}), goal directed (\mathcal{F}^{goal}), and group (\mathcal{F}^{group}) forces influencing the individuals.

ownership and social relations described in Section 4. The mechanism adopted is sufficiently general to accommodate external information (e.g. the state of alert, time of day) alongside information on the observed scene in determining whether or not to raise an alarm.

Three increasingly sophisticated definitions are considered for what constitutes an abandoned bag. The first adopts the simple *baseline* definition that defined the *PETS (2006)* challenge. In this, a threat (i.e. abandonment) is defined as follows:

- Bag unattended if no person within 2 m.
- Bag abandoned if unattended for 30 s.

Here, the notions of ownership and social relationships are not used.

The second definition (*owner*) includes the notion of ownership (Section 4.1) and is defined as follows:

- Bag unattended if owner is not within 2 m.
- Bag unattended if there is no assigned owner and if no person within 2 m.
- Bag abandoned if unattended for 30 s.

When there is no assigned owner, this is equivalent to the baseline definition, but where one or more possible owners have been assigned, the condition for an alarm to be raised is less stringent since the behaviours of non-owners within the scene is ignored.

The third definition (*owner + group*) includes both the notions of ownership (Section 4.1) and social relationships (Section 4.2). In this, a threat is defined as follows:

- Bag unattended if owner and social group of owner are not within 2 m.
- Bag unattended if there is no assigned owner and if no person within 2 m.
- Bag abandoned if unattended for 30 s.

This relaxes the *owner* definition in the direction of the *baseline* definition, since now the circle of people attending to a bag is widened to include people in the same group as the possible owner (s). The likelihood of raising an alarm is therefore reduced.

5.1. Implementation

The aim in threat assessment is to make it straightforward to encode the evolving state of the world and explore different behav-

oural patterns that constitute a potential threat. To achieve this, a simple logic-based inference system (Prolog) is adopted in which the current state of the world is represented by a set of facts and the behavioural patterns that constitute potential threats are encoded as rules.

The elements of this logic-based approach are:

- Facts (logical atoms), which are employed to describe situations. A fact is of the form $R(A, B, \dots)$, where R indicates a type of relation between the elements inside the brackets.
- Rules, which are employed to infer new facts from existing ones.

Given these elements, the threat assessment proceeds in two steps:

1. Tracking and detection data are converted into a set of facts.
2. A set of pre-defined rules is invoked to infer additional facts.

The position of an object in each frame is represented by a unique ID for the object, its class (person or bag), its x, y position on the ground-plane and the frame number:

$track(id, class, x, y, frame)$.

The social relationships between individuals are represented by a single predicate that records a unique group ID for each person. This partitions the set of people into social groups. Any person not assigned to a social group is assumed to be outside any group. This is represented simply by facts of the form:

$group(id, group_id)$.

For convenience, a ‘class’ predicate is used (as in $class(id, person)$) to record the class of each object independently of the ‘track’ facts.

The ownership of bags is inferred next by a set of Prolog rules that embody the criteria described in above. The result is a new set of facts, each representing the ownership of a bag (b) by a person (p):

$owner(p, b)$.

Finally, the alarm condition for the chosen threat definition is posed as a Prolog query. As part of this, for the baseline definition, the condition that a bag is attended translates into the rule:

$attended(B, T) :- class(P, person), nearby(P, T, B, T, 2)$.

Here the rule states that a bag is attended at time T (shown on the left of the ‘:-’) if there is a person (call them P), and the position of P at time T is within 2 m (i.e. nearby) of the position of B at time T (shown on the right of the ‘:-’). Upper case arguments are used to signify that these are variables.

The equivalent set of rules for the *owner + group* definition, incorporating the notions of ownership and social relationships, is as follows:

```
attended(B, T) :- owner(P, B), nearby(P, T, B, T, 2), !.
attended(B, T) :-
  \ + owner(., B), track(P, person, ., ., T), nearby(P, T, B, T, 2).
attended(B, T) :- owner(P, B), knows(P, Q), nearby(Q, T, B, T, 2), !.
knows(P, Q) :- group(P, G), group(Q, G).
```

The first rule states that a bag B is attended at time T if there is an owner P for the bag and this person is nearby. The second rule invokes the baseline notion of being attended when there is no owner - the meaning of ‘+’ before the owner predicate means that this is not present in the database. The third rule states that a bag is attended (at time T) if there is a second person Q who is nearby the bag and P and Q know one another. The fourth rule implements the notion of two people knowing one another in terms of their group membership - i.e. they know one another if they are from the same social group. The *owner* definition, incorporating only the notion of ownership, is defined by the first two of the rules above.

Finally the condition for an alarm to be raised is the same for all three definitions - a bag must be unattended for a fixed period of time. The definition of ‘unattended’ is expressed in terms of the different definitions of attended, as follows:

```
unattended(B, T) :- class(B, bag), track(B, bag, ., ., T),
  \ + attended(B, T).
```

This states that an object is unattended at time T if it is a bag, it is in existence at time T , and there is no ‘attended’ fact in the database for that bag at time T .

Thus, only the definition of ‘attended’ varies between the three definitions of what constitutes an alarm.

Generally, Prolog was found to be a convenient way to represent definitions in a readily understood fashion, facilitating extension and experimentation. On the other hand, there are aspects of the inference mechanism in Prolog that require care - for example the use of the cut (!) in two of the rules above is necessary to avoid the same alarm being raised multiple times.

6. Results

6.1. Datasets

Two different datasets are used to test the performance of the proposed algorithms (see Fig. 5), the publicly available PETS2006 (PETS, 2006) and the second produced during the SUBITO project specifically for this study. The PETS2006 dataset consists of ten sequences with increasing complexity of a staged abandoned bag scenario at a train station. All four camera views in the dataset were used in turn for the first four sequences used (PETS-S1-1, PETS-S1-2, PETS-S1-3 and PETS-S1-4), and camera view 3 used only for the other sequences (PETS-S2-3, PETS-S3-3, PETS-S4-3, PETS-S5-3, PETS-S6-3 and PETS-S7-3). The SUBITO dataset was recorded specifically for the SUBITO project. It contains thirteen sequences (19–22, 24–29, 31, 36, 37) each recorded from four synchronised cameras placed around the scene. In sequences 19–22 a single person brings a bag to a marked position and loiters around the bag (sequence 19), abandons the bag (sequence 20), or leaves the bag unattended for a while and then comes back (sequences 21, 22). Sequences 24–29, 31, 36 and 37 contain more challenging variants in terms of number of people and the group relationships. Each action is recorded 12 times for different entrance/exit directions. Depending on different threat definitions, the same action may or may not raise an alarm. Each sequence therefore should either correspond to 12 alarms (except for sequence 36 which only corresponds to 11 alarms), or none. The ground-truth alarms were obtained manually for all three threat definitions. The alarm time is determined by first visually deciding the very frame when the owner is just outside the prescribed distance from the bag, then adding a fixed time interval before the alarm is raised. Within the SUBITO dataset, the critical distance around a bag is assumed to be 2.5 m (as opposed to 2 m used in the PETS2006 challenge)- this assumption is therefore used in the three threat definitions. The time a bag must remain unattended to raise an alarm is reduced to 4 s.

6.2. Preliminary experiments on PETS2006 data

In the first experiments, the baseline functionality of (PETS, 2006) was implemented and evaluated. These experiments were carried out using an earlier version of the threat assessment logic implemented in C++. This was subsequently re-implemented in Prolog as part of the real-time system. To achieve this, the Prolog is queried for an alarm on every frame, based on the current state of the world and pertinent facts from the recent past. This world

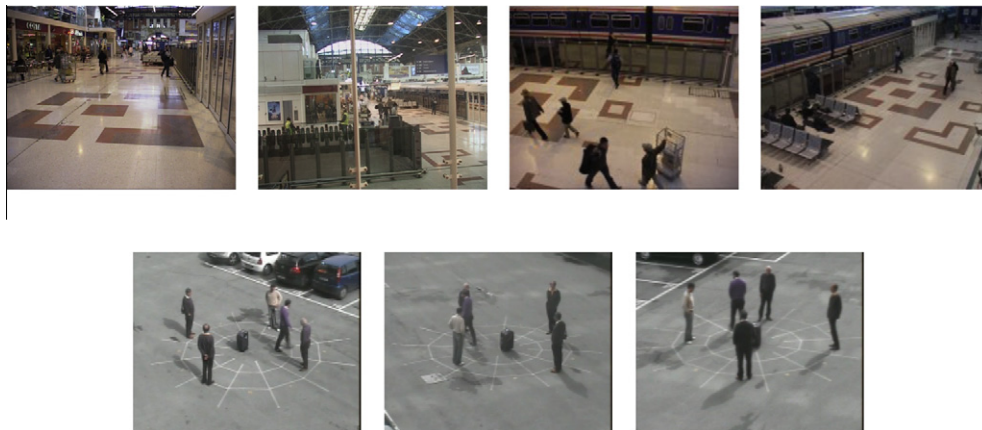


Fig. 5. Datasets used. Top row: four views from PETS2006 which contains scenarios with abandoned luggage. Bottom row: three views from the SUBITO dataset describes scenarios where luggage owner enters the scene, sometimes interacts with other individuals and leaves the scene with/without the luggage.

Table 1

Aggregate results across all SUBITO sequences comparing predicted alarms with corresponding baseline/owner/group ground truth.

Ruleset	TP	GTalarms	Alarms	Recall	Precision
baseline	16	71	35	0.23	0.46
owner	48	143	75	0.34	0.64
owner+group	39	107	66	0.36	0.59

Table 2

Aggregate results across all SUBITO sequences comparing the use of all three threat definitions with the ground truth for the *owner + group* definition.

Ruleset	TP	GTalarms	Alarms	Recall	Precision
baseline	16	107	35	0.15	0.46
owner	42	107	75	0.39	0.56
owner+group	39	107	66	0.36	0.59

model is continually refreshed with the current location of each tracked object.

For the threat assessment to be correct, the system is required to raise an alarm following a potential threat, and to correctly identify the ID of the abandoned bag. Specifically, an alarm must be raised within 50 frames of a ground-truth alarm for it to be successful detected. The results on the PETS2006 dataset employ automatic tracking using an implementation of Breitenstein et al. (2011) and bag detection using Porikli et al. (2008). Alarms were raised correctly on all tested sequences except PETS-S4-3 and PETS-S7-3. The failures on these two sequences were caused by individuals, having nothing to do with the abandoned bag, nevertheless being close enough to prevent the bag being classified as unattended. This result motivates the concept of ownership considered in the main set of experiments.

6.3. Experiments on SUBITO data

The main set of experiments were carried out on the challenging SUBITO dataset. The inverse SFM system is run in batch mode so that it has access to an entire sequence in predicting social groups rather than only the history up until the current time. The entire sequence is therefore used in inferring the set of alarms. This enabled evaluation of the interaction of the detection and tracking sub-system and the threat assessment sub-system, giving the inverted SFM the best chance of assigning correct social groups within relatively short scenarios. A single threshold in the inverse SFM system controls the propensity of pairs of individuals to be combined into the same group; a lower threshold results in larger social groups. For the SUBITO data, we found that both precision and recall reach their highest values within a small range of this threshold and the results we present are for a choice of threshold in this range.

The aggregate results across all SUBITO sequences are shown in Table 1, comparing predicted alarms with the corresponding ground-truth - that is baseline results are compared with the baseline ground-truth, etc. The aggregate results comparing the use of all three threat definitions with the ground-truth for the *owner + group* definition are shown in Table 2. As expected, the precision and recall for the *baseline* definition are lower in this case since the ground-truth reflects a more sophisticated notion of threat, incorporating concepts that are not present in the *baseline* definition. The evaluation reported here attended only to the time an alarm is raised and ignored the ID for the person and bag involved. Where there is more than one true positive alarm for a ground-truth alarm, this is counted once in computing recall and does not contribute to loss of precision. In other words, multiple

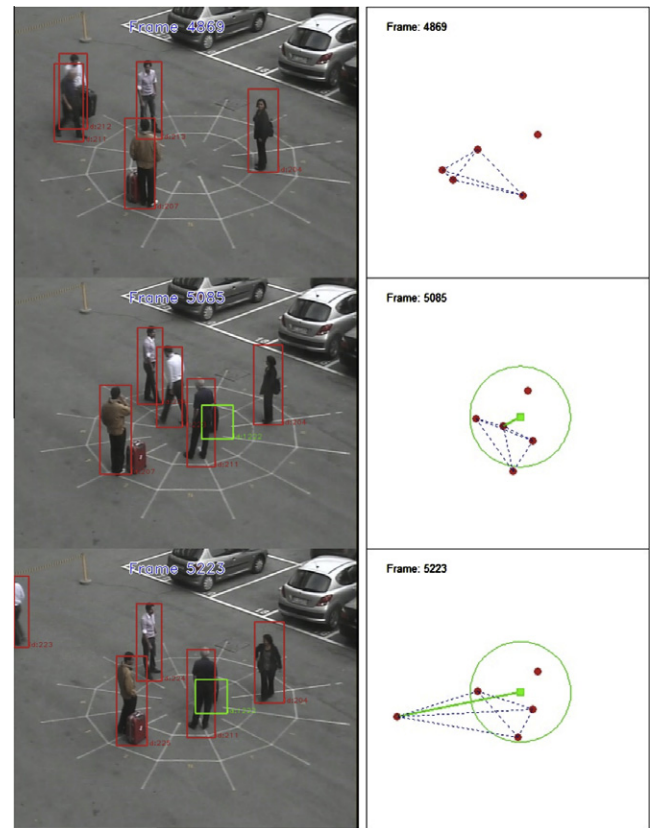


Fig. 6. Social group analysis applied to SUBITO sequence 36 resulting in correct suppression of false alarm.

predicted alarms for the same ground-truth alarm are counted only once. In general, there were few instances of this occurring in the experiment.

Within Table 2, there is a clear improvement in precision and recall between *baseline* and *owner* definitions. However, the comparison of performance between *owner* and *owner + group* definitions is less decisive. Here the recall has reduced slightly with the introduction of the social relationships, but there is a comparable improvement in precision. Looking in more detail at the results on individual sequences and alarms, several alarms have been suppressed by correct assignment of an owner and partner to the same social group. This is illustrated in Fig. 6 showing a set of frames from SUBITO sequence 36. Two individuals (d:211, d:212) entering the scene (Fig. 6 (top)) are assigned to the same social group (indicated by blue line between them), and one is detected as the owner of a bag (d:212) that appears within the scene (Fig. 6 (middle)). The owner subsequently goes away from the bag and outside the prescribed distance (shown as a green circle around the bag), leaving their partners attending to the bag (Fig. 6 (bottom)). No alarm is raised.

In general the recall and precision are below acceptable performance for a deployed threat assessment system. The principal source of error arises from the highly challenging video sequences containing multiple overlapping actors at any time. The consequential limitations in detection and tracking performance are translated directly into the threat assessments that can be achieved using the logic described above. Some improvement in performance was achieved by automatically stitching together tracks for which there is sufficient evidence that they belong to the same objects at different periods of time - specifically, one track (of more than 10 frames duration) ends within 4 s and 1 m of another track (of more than 10 frames duration) beginning.

Table 3

Aggregate results across all SUBITO sequences comparing the use of all three threat definitions with the ground truth for the *owner+group* definition with *stiched-together tracks*.

Ruleset	TP	GTalarms	Alarms	Recall	Precision
baseline	15	107	36	0.14	0.42
owner	43	107	94	0.40	0.46
owner+group	41	107	88	0.38	0.47

The precision and recall for the equivalent evaluation to that in Table 2 is shown in Table 3. Finally, a real-time system that incorporates all stages of the pipeline, including on-line estimation of social groups up to the current frame, has also been implemented to demonstrate the practical viability of the method.

7. Conclusions and future work

This paper has described a video surveillance framework that detects abandoned objects in surveillance scenes containing multiple interacting individuals, extending the state of the art. Future work will address methods to further improve the underpinning object (person and bag) detection and tracking accuracy, as well as introduction of goal-directed and intentionality modelling strategies in the behavioural analysis.

There is scope to perform a more rigorous analysis of ownership through detecting bags being carried into the scene and hence identifying the owner more reliably. Similarly, confidence that a bag has been removed from the scene would be raised if it could be detected as it was carried out. There is prior work on this problem that should in principle be directly applicable to sequences such as those in the SUBITO dataset (e.g. Damen and Hogg (2008)).

Finally, expressing the conditions of a threat in terms of logic, suggests that it may be possible to induce such conditions automatically from examples, thereby providing a way to incorporate different kinds of information about the scene without having to provide the logical rules by hand. Earlier work on the use of inductive logic programming in video analysis indicates how this might be achieved in principle (Dubba et al., 2010).

Acknowledgements

This work was supported by EC projects SUBITO Grant Agreement No. 218004 and FP7-ICT-247022 MASH. Any opinions expressed in this paper do not necessarily reflect the views of the European Community. The Community is not liable for any use that may be made of the information contained herein.

References

Ardö, H., Aström, K., 2007. Multi sensor loitering detection using online viterbi. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1423-9.

Arsić, D., Hofmann, M., Schuller, B., Rigoll, G., 2007. Multi-camera person tracking and left luggage detection applying homographic transformation. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1423-9.

Auvinet, E., Grossmann, E., Rougier, C., Dahmane, M., Meunier, J., 2006. Left-luggage detection using homographies and simple heuristics. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.

Aveni, A.F., 1977. The not-so-lonely crowd: friendship groups in collective behavior. *Sociometry* 40 (1), 96–99.

Bayona, A., SanMiguel, J.C., Martínez, J.M., 2009. Comparative evaluation of stationary foreground object detection algorithms based on background subtraction techniques. In Proc. of the 6th IEEE Internat. Conference on Advanced Video and Signal Based Surveillance (AVSS 09), pp. 25–30. <http://dx.doi.org/10.1109/AVSS.2009.35>.

Berclaz, J., Shahroki, A., Fleuret, F., Ferryman, J., Fua, P., 2009. Evaluation of probabilistic occupancy map people detection for surveillance systems. In Proc. of the IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS), pp 55-62. ISBN 978-07049-1501-4.

Birch, P., Hassan, W., Bangalore, N., Young, R., Chatwin, C., 2011. Stationary traffic monitor. In Proc. 4th Internat. Conf. on Imaging for Crime Detection and Prevention (ICDP-11), pp. 1–6. <http://dx.doi.org/10.1049/ic.2011.0128>.

Blackman, S.S., 2004. Multiple hypothesis tracking for multiple target tracking. *IEEE Aerospace and Electronic Systems Magazine* 19 (1), 5–18.

Brand, M., Oliver, N., Pentland, A., 1997. Coupled hidden markov models for complex action recognition. In Proc. Computer Vision and Pattern Recognition, pp. 994–999.

Breitenstein, M.D., Reichlin, F., Leibe, B., Koller-Meier, E., Van Gool, L., 2011. Online multiperson tracking-by-detection from a single, uncalibrated camera. *IEEE Trans. on Pattern Analysis and Machine Intelligence* 33 (9), 820–833.

Challenger, R., Clegg, C.W., Robinson, M.A., Leigh, M., 2009. Understanding Crowd Behaviours: Simulation Tools. Technical Report, University of Leeds.

Dalley, G., Wang, X., Grimson, W.E.L., 2007. Event detection using an attention-based tracker. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1423-9.

Damen, D., Hogg, D., 2008. Detecting carried bags in short video sequences. In Proc. 10th European Conf. on Computer Vision, vol. 5304, pp. 154–167.

Damen, D., Hogg, D., 2012. Explaining activities as consistent groups of events: A bayesian framework using attribute multiset grammars. *International Journal of Computer Vision* 98 (1), 83–102. <http://dx.doi.org/10.1007/s11263-011-0497-0>.

Dubba, K.S.R., Cohn, A.G., Hogg, D.C., 2010. Even model learning from complex videos using ILP. In Proceedings ECAI, vol. 215, pp. 93–98.

Evangelio, R., Sikora, T., 2011. Static object detection based on a dual background model and a finite-state machine. *EURASIP Journal on Image and Video Processing*, <http://jivp.erasipjournals.com/content/2011/1/858502>.

Ge, W., Collins, R., Ruback, B., 2009. Automatically detecting the small group structure of a crowd. In Workshop on Applications of Computer Vision, pp. 1–8.

Guler, S., Farrow, M.K., 2006. Abandoned object detection in crowded places. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.

Helbing, D., Molnar, P., 1995. Social Force Model for Pedestrian Dynamics. *Physical Review E. Statistical, Nonlinear, and Soft Matter Physics* 51, 4282.

Hongeng, S., Nevatia, R., Brémond, F., 2004. Video-based event recognition: activity representation and probabilistic recognition methods. *Computer Vision and Image Understanding* 96, 129–162.

i-LIDS Imagery Library for Intelligent Detection Systems, <http://www.ilids.co.uk>. Last accessed: 29 January 2012.

Ivanov, Y., Bobick, A., 2000. Recognition of visual activities and interactions by stochastic parsing. *IEEE Trans. on Pattern Analysis and Machine Intelligence* 22 (8), 852–872.

Jacques, J., Braun, A., Soldera, J., Musse, S., Jung, C., 2007. Understanding people motion in video sequences using voronoi diagrams. *Pattern Analysis and Applications* 10, 321–332.

Joo, S.-W., Chellappa, R., 2006. Attribute Grammar-Based Event Recognition and Anomaly Detection. In Proceedings Internat. Workshop on Semantic Learning Applications in Multimedia, New York, NY June.

Krahnstoever, N., Tu, P., Sebastian, T., Perera, A., Collins, R., 2006. Multi-view detection and tracking of travelers and luggage in mass transit environments. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.

Leung, V., Herbin, S., 2011. Flexible tracklet association for complex scenarios using a markov logic network. In Proc. 11th Internat. Workshop on Visual Surveillance, pp. 1870–1875.

Li, L., Luo, R., Ma, R., Huang, W., Leman, K., 2006. Evaluation of an IVS system for abandoned object detection on PETS 2006 datasets. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.

Li, J., Gong, S., Xiang, T., 2008. Global behaviour inference using probabilistic latent semantic analysis. In Proc. British Machine Vision Conf., DOI:10.5244/C.22.20.

Lu, S., Zhang, J., Feng, D., 2007. Detecting unattended packages through human activity recognition and object association. *Pattern Recognition* 8, 2173–2184.

Lv, F., Song, X., Wu, B., Singh, V.K., Nevatia, R., 2006. Left-luggage detection using bayesian inference. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.

Martínez-del-Rincón, J., Herrero-Jaraba, J., Ral Gómez, J., Orrite-Uruñuela, C., 2006. Automatic left luggage detection and tracking using multi-camera UKF. In Proceedings IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.

Moussaid, M., Helbing, D., Garnier, S., Johansson, A., Combe, M., Theraulaz, G., 2009. Experimental study of the behavioural mechanisms underlying self-organization in human crowds. *Proc. of the Royal Society B: Biological Sciences* 276 (1668), 2755–2762.

Moussaid, M., Perozo, N., Garnier, S., Helbing, D., Theraulaz, G., 2010. The walking behaviour of pedestrian social groups and its impact on crowd dynamics. *PLoS One* 5 (4), e10047.

Pellegrini, S., Ess, A., Van Gool, L., 2010. Improving data association by joint modeling of pedestrian trajectories and groupings. *Proc. European Conference on Computer Vision*, vol. 6311. Springer, Berlin/Heidelberg, pp. 452–465.

PETS. 2006. In Ninth IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS), New York, USA, June 18. ISBN 0-7049-1422-0.

PETS. 2007. In Tenth IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS), Rio De Janeiro, October 14. ISBN 0-7049-1423-9.

- Porikli, F., Yin, Z., 2006. Temporally static region detection in multi-camera systems. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.
- Porikli, F., Ivanov, Y., Haga, T., 2008. Robust abandoned object detection using dual foregrounds. EURASIP Journal on Advances in Signal Processing. Article ID 197875. Available: <http://asp.eurasipjournals.com/content/2008/1/197875>.
- Reynolds, C.W., 1987. Flocks, herds and schools: A distributed behavioral model. In Proc. of the 14th Annual Conference on Computer Graphics and Interactive Techniques, pp. 25–34. <http://dx.doi.org/10.1145/37401.37406>.
- Ribeiro, P.C., Moreno, P., Santos-Victor, J., 2007. Detecting luggage related behaviors using a new temporal boost algorithm. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1423-9.
- Richardson, M., Domingos, P., 2006. Markov logic networks. Machine Learning 62, 107–136.
- Rudloff, C., Matyus, T., See, S., Bauer, D., 2011. Can walking behaviour be predicted? an analysis of the calibration and fit of pedestrian models. In 90th Annual Meeting of the Transportation Research Board, January.
- Shet, V.D., Harwood, D., David, L.S., 2005. Vidmap: Video monitoring of activity with prolog. In Proc. IEEE Conference on Advance Video and Signal Based Surveillance, pp. 224–229.
- Shi, Y., Huang, Y., Minnen, D., Bobick, A., Essa, I., 2004. Propagation networks for recognition of partially ordered sequential action. Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition 2, 862–869. <http://dx.doi.org/10.1109/cvpr.2004.1315255>.
- Singh, H., Arter, R., Dodd, L., Langston, P., Lester, E., Drury, J., 2009. Modelling subgroup behaviour in crowd dynamics dem simulation. Applied Mathematical Modelling 33 (12), 4408–4423.
- Smith, K., Quelhas, P., Gatica-Perez, D., 2006. Detecting abandoned luggage items in a public space. In Proc. IEEE Internat. Workshop on Performance Evaluation of Tracking and Surveillance (PETS). ISBN 0-7049-1422-0.
- Sochman, J., Hogg, D.C., 2011. Who knows who - inverting the social force model for finding groups. In Proc. IEEE Intelligent Workshop on Socially Intelligent Surveillance and Monitoring.
- Tian, Y., Feris, R., Liu, H., Humpapur, A., Sun, M.-T., 2010. Robust detection of abandoned and removed objects in complex surveillance videos. In Proc. IEEE Trans. on Systems, Man, and Cybernetics, Part C: Applications and Reviews, pp. 1–12.
- Tran, S., Davis, L., 2008. Visual event modelling and recognition using markov logic networks. In Proc. European Conf. on Computer Vision, pp. 610–623.
- Wang, X., Ma, X., Grimson, E., 2009. Unsupervised activity perception in crowded and complicated scenes using hierarchical bayesian models. IEEE Trans. on Pattern Analysis and Machine Intelligence 31 (3), 539–555.
- Yildiz, A., Akgul, Y.S., 2010. A fast method for tracking people with multiple cameras, Technical Report, Vision Lab, Gebze Institute of Technology.
- Zivkovic, Z., 2004. Improved adaptive gaussian mixture model for background subtraction. Proc. Internat. Conf. on Pattern Recognition 2, 28–31.