



**HAL**  
open science

# Video Understanding Framework For Automatic Behavior Recognition

François Bremond, Monique Thonnat, Marcos Zuniga

► **To cite this version:**

François Bremond, Monique Thonnat, Marcos Zuniga. Video Understanding Framework For Automatic Behavior Recognition. Behavior Research Methods, 2006, 3 (38), pp.416-426. inria-00276938

**HAL Id: inria-00276938**

**<https://inria.hal.science/inria-00276938v1>**

Submitted on 2 May 2008

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Video Understanding Framework For Automatic Behavior Recognition

François Brémond, Monique Thonnat, Marcos Zúñiga  
INRIA Sophia Antipolis, ORION group  
2004, route des Lucioles, BP93  
06902 Sophia Antipolis Cedex — France  
(33) 4 9238 7657  
firstname.surname@sophia.inria.fr

We propose an activity monitoring framework based on a platform called VSIP, enabling behavior recognition in different environments. To allow end-users to actively participate in the development of a new application, VSIP separates algorithms from a priori knowledge. For describing how VSIP works, we present a full description of a system developed with this platform for recognizing behaviors, involving either isolated individual, group of people or crowds, in the context of visual monitoring of metro scenes using multiple cameras. In this work, we also illustrate the capability of the framework to easily combine and tune various recognition methods dedicated to the visual analysis of specific situations (e.g. mono/multi actors activities, numerical/symbolic actions or temporal scenarios). We also present other applications using this framework, in the context of behavior recognition. VSIP has shown a good performance on human behavior recognition for different problems and configurations, being suitable to fulfill a large variety of requirements.

One of the most challenging problems in the domain of computer vision and artificial intelligence is video understanding. The research in this area concentrates mainly on the development of methods for analysis of visual data in order to extract and process information about the behavior of physical objects in a scene.

Most approaches in the field of video understanding incorporated methods for detection of domain-specific events. Examples of such systems use Dynamic Time Warping for gesture recognition (Bobick & Wilson, 1997) or self-organizing networks for trajectory classification (Owens & Hunter, 2000). The main drawback of these approaches is the usage of techniques specific only to a certain domain which causes difficulties on applying these techniques to other areas. Therefore some researchers have adopted a two-steps approach to the problem of video understanding:

1. A visual module is used in order to extract visual cues and primitive events.
2. This information is used in a second stage for the detection of more complex and abstract behavior patterns (Hu, Tan, Wang, & Maybank, 2004).

By dividing the problem into two sub-problems we can use simpler and more domain-independent techniques in each step. The first step makes usually extensive usage of stochastic methods for data analysis while the second step conducts structural analysis of the symbolic data gathered at the preceding step (see Figure 1). Examples of this two-level architecture can be found in the works of (Ivanov & Bobick, 2000) and (Vu, Brémond, & Thonnat, 2003).

This approach is available as a platform for image sequence understanding called VSIP (Video Surveillance Interpretation Platform) which was developed at the research group ORION at INRIA (Institut National de Recherche en Infor-

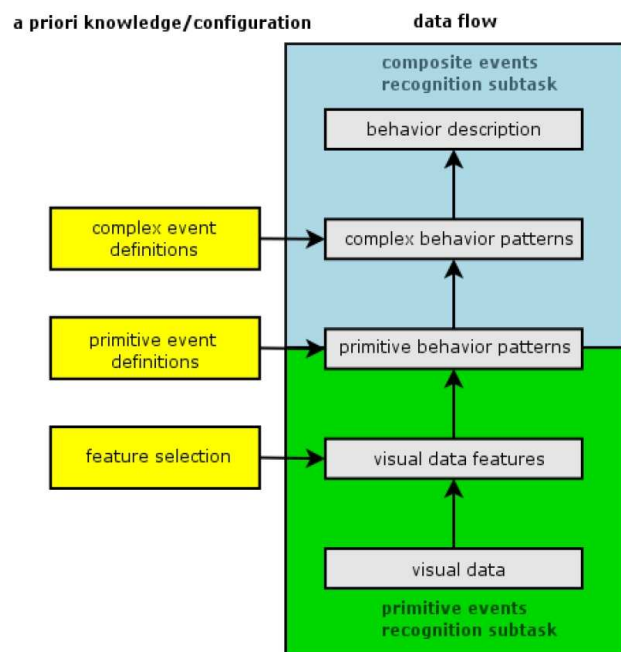


Figure 1. A general architecture of a video understanding system. The steps depicted in the figure describe the data flow during a video understanding process.

matique et en Automatique), Sophia Antipolis. VSIP is a generic environment for combining algorithms for processing and analysis of videos which allows to flexibly combine and exchange various techniques at the different stages of the

video understanding process. Moreover, VSIP is oriented to help developers describing their own scenarios and building systems capable of monitoring behaviors, dedicated to specific applications.

At the first level, VSIP extracts primitive geometric features like areas of motion. Based on them, objects are recognized and tracked. At the second level those events in which the detected objects participate, are recognized. For performing this task, a special representation of events is used which is called event description language (Vu et al., 2003). This formalism is based on an ontology for video events presented in (Brémond, Maillot, Thonnat, & Vu, 2004) which defines concepts and relations between these concepts in the domain of human activity monitoring. The major concepts encompass different object types and the understanding of their behavior from the point of view of the domain expert.

In this article, we illustrate this framework by presenting a system developed using VSIP platform for recognizing people behaviors, such as fighting or vandalism, occurring in a metro scene, viewed by one or several cameras. This work has been performed in the framework of the European project ADVISOR (<http://www-sop.inria.fr/orion/ADVISOR>). Our goal is to recognize in real time behaviors involving either isolated individuals, groups of people or crowds from real world video streams coming from metro stations. To reach this goal, we have developed a system which takes as input video streams coming from cameras and generates annotation about the behaviors recognized in the video streams.

This article is organized as follows. In the second section, we will present the state of the art on behavior recognition. In the third section, we will briefly present the global system and its vision module. Then details of the behavior recognition process will be illustrated with five behavior recognition examples—*fighting*, *blocking*, *vandalism*, *overcrowding*, and *fraud* behaviors—and with an analysis of the obtained results. In the fourth section, we will discuss VSIP’s capability for dealing with other applications under various configurations.

## RELATED WORK

Since the years 90s, a problem of focus in cognitive vision has been Automatic Video Interpretation. There are now several research units and companies defining new approaches to design systems that can understand specific scenarios in dynamic scenes. We define a scenario as a combination of states, events or sub scenarios. Behaviors are specific scenarios, defined by the users.

Three main categories of approaches are used to recognize scenarios:

1. Probabilistic/neural network combining potentially recognized scenario.
2. Symbolic network that Stores Totally Recognized Scenarios.
3. Symbolic network that Stores Partially Recognized Scenarios.

For the computer vision community, a natural approach consists in using a probabilistic/neural network. The nodes of

this network correspond usually to scenarios that are recognized at a given instant with a computed likelihood. For example, (Howell & Buxton, 2002) proposed an approach to recognize a scenario based on a neural network (time delay Radial Basis Function). (Hongeng, Brémond, & Nevatia, 2000) proposed a scenario recognition method that uses concurrence Bayesian threads to estimate the probability of potential scenarios.

For the artificial intelligence community, a natural way to recognize a scenario is to use a symbolic network which nodes correspond usually to the boolean recognition of scenarios. For example, (Rota & Thonnat, 2000) used a declarative representation of scenarios defined as a set of spatio-temporal and logical constraints. They used a traditional constraint resolution technique to recognize them. To reduce the processing time for the recognition step, they proposed to check the consistency of the constraint network using the AC4 algorithm. (Gerber, Nagel, & Schreiber, 2002) defined a method to recognize a scenario based on a fuzzy temporal logic. The common characteristic of these approaches is that all totally recognized behaviors are stored (recognized in the past) (Vu et al., 2003).

Another approach consists in using a symbolic network and storing partially recognized scenarios (to be recognized in the future). For example, (Ghallab, 1996) has used the terminology chronicle to express a temporal scenario. A chronicle is represented as a set of temporal constraints on time-stamped events. The recognition algorithm keeps and updates partial recognition of scenarios using the propagation of temporal constraints based on RETE algorithm. Their applications are dedicated to the control of turbines and telephonic networks. (Chleq & Thonnat, 1996) made an adaptation of temporal constraints propagation for video surveillance. In the same period, (Pinhanez & Bobick, 1998) have used Allen’s interval algebra to represent scenarios and have presented a specific algorithm to reduce its complexity.

All these techniques allow an efficient recognition of scenarios, but there are still some temporal constraints which cannot be processed. For example, most of these approaches require that the scenarios are bounded in time (Ghallab, 1996), or process temporal constraints and atemporal constraints in the same way (Rota & Thonnat, 2000).

Another problem that has captured the attention of researchers recently is the problem of unsupervised behavior learning and recognition, consisting in the capability of a vision interpretation system of learning and detecting the frequent scenarios of a scene without requiring the prior definition of behaviors by the user. The unsupervised behavior learning and recognition problem in the field of computer vision is addressed only in a few works. Most of the approaches are confined to a specific domain and take advantage of domain knowledge in order, for example, to choose a proper model or to select features. One of the most widely used techniques for learning scenarios in an unsupervised manner is the topology of a Markov model. (Brand & Kettner, 2000) use an entropy-based function instead of the Maximum-Likelihood estimator in the E-step of the EM-algorithm for learning parameters of Hidden Markov Mod-

els (HMM). This leads to a concentration of the transitional probabilities just on several states which correspond in most of the cases to meaningful events. Another approach is based on variable length Markov models which can express the dependence of a Markov state on more than one previous state (Galata, Cohn, Magee, & Hogg, 2002). While this method learns good stochastic models of the data it cannot handle temporal relations. A further similar technique is based on hierarchical HMMs whose topology is learned by merging and splitting states (Xie, Chang, Divakaran, & Sun, 2003). The advantage of the above techniques for topology learning of Markov models is that they work in a completely unsupervised way. Additionally, they can be used after the learning phase to recognize efficiently the discovered events. On the other hand, these methods deal with simple events, are not capable of creating concept hierarchies and there is no guaranty that the states of these models correspond to meaningful events.

A different approach for this problem was proposed by (Toshev, Brémond, & Thonnat, 2006). In this work, the A priori algorithm from the field of data mining is used to propose a method for unsupervised learning of behaviors from videos. The developed algorithm processes a set of generic primitive events and outputs the frequent patterns of these primitive events, also interpreted as frequent composite events. In a second step, models of these composite events are automatically generated (i.e. learned) in the event description language defined by (Vu et al., 2003), which can be used to recognize the detected composite events in new videos. This application was used for detecting frequent behaviors on a parking lot monitoring system.

This review of the state of the art shows the large diversity of video understanding techniques in automatic behavior recognition. The challenge is to combine efficiently these techniques to address the large diversity of the real world.

## VIDEO UNDERSTANDING PLATFORM

The video interpretation platform is based on the cooperation of a vision and a behavior recognition module as shown on Figure 2.

The vision module is composed of three tasks. First a motion detector and a frame to frame tracker generates a graph of mobile objects for each calibrated camera. Second, a combination operation is performed to combine the graphs computed for each camera into a global one. Third, this global graph is used for long term tracking of individuals, vehicles, groups of people and crowds as the scene evolves.

For each tracked actor, the behavior recognition module performs three levels of reasoning: states, events and scenarios. On top of that, we use 3D scene models (i.e. geometric model of the empty scene, including the furniture), one for each camera, as a priori contextual knowledge of the observed scene. We define in a scene model the 3D positions and dimensions of the static scene objects (e.g. a bench, a ticket vending machine) and the zones of interest (e.g. an entrance zone). Semantic attributes (e.g. fragile) can be associated

to the objects or zones of interest to be used in the behavior recognition process.

On this paper we focus on the behavior recognition process, as it is our current focus of interest<sup>1</sup>. The goal of this process is to recognize specific behaviors occurring in an observed scene. A main problem in behavior recognition is the ability to define and reuse methods to recognize specific behaviors, knowing that the perception of behaviors is strongly dependent on the site, the camera view point and the individuals involved in the behaviors. Our approach consists in defining a formalism allowing us to write and easily reuse all methods needed for the recognition of behaviors. This formalism is based on three main ideas:

1. The formalism should be flexible enough to allow various types of *operators* to be defined (e.g. a temporal filter or an automaton). We use *operator* as an abstract term to define programs. This term will be defined in the following section.
2. All the needed knowledge for an operator should be explained within the operator so that it can be easily reused.
3. The description of the operators should be declarative in order to build an extensible library of operators.

### *Behavior representation*

We call an actor of a behavior any scene object involved in the behavior, including static objects (equipment, zones of interest), individuals, groups of people or crowds. The entities needed to recognize behaviors correspond to different types of concepts which are:

1. **The basic properties:** A characteristic of an actor such as its trajectory or its speed.
2. **The states:** A state describes a situation characterizing one or several actors defined at time  $t$  (e.g. a group is agitated) or a stable situation defined over a time interval. For the state: “an individual stays close to the ticket vending machine”, two actors are involved: an individual and a piece of equipment.
3. **The events:** An event is a change of states at two consecutive times (e.g. a group enters a zone of interest).
4. **The scenarios:** A scenario is a combination of states, events or sub scenarios. Behaviors are specific scenarios (dependent on the application) defined by the users. For example, to monitor metro stations, end-users have defined five targeted behaviors: “Fraud”, “Fighting”, “Blocking”, “Vandalism” and “Overcrowding”.

To compute all the needed entities for the recognition of behaviors, we use a generic framework based on the definition of *operators* which are program descriptions containing four elements:

1. **Name:** Indicates the entity to be computed such as the state “an individual is walking” or “the trajectory is straight”.
2. **Input:** Gives a description of input data. There are two types of input data: basic properties characterizing an actor and sub entities computed by other *operators*.
3. **Body:** Contains a set of competitive methods to compute the entity. All these methods are able to compute this

<sup>1</sup> For more details on the vision processing module, see (Cupillard, Brémond, & Thonnat, 2004).

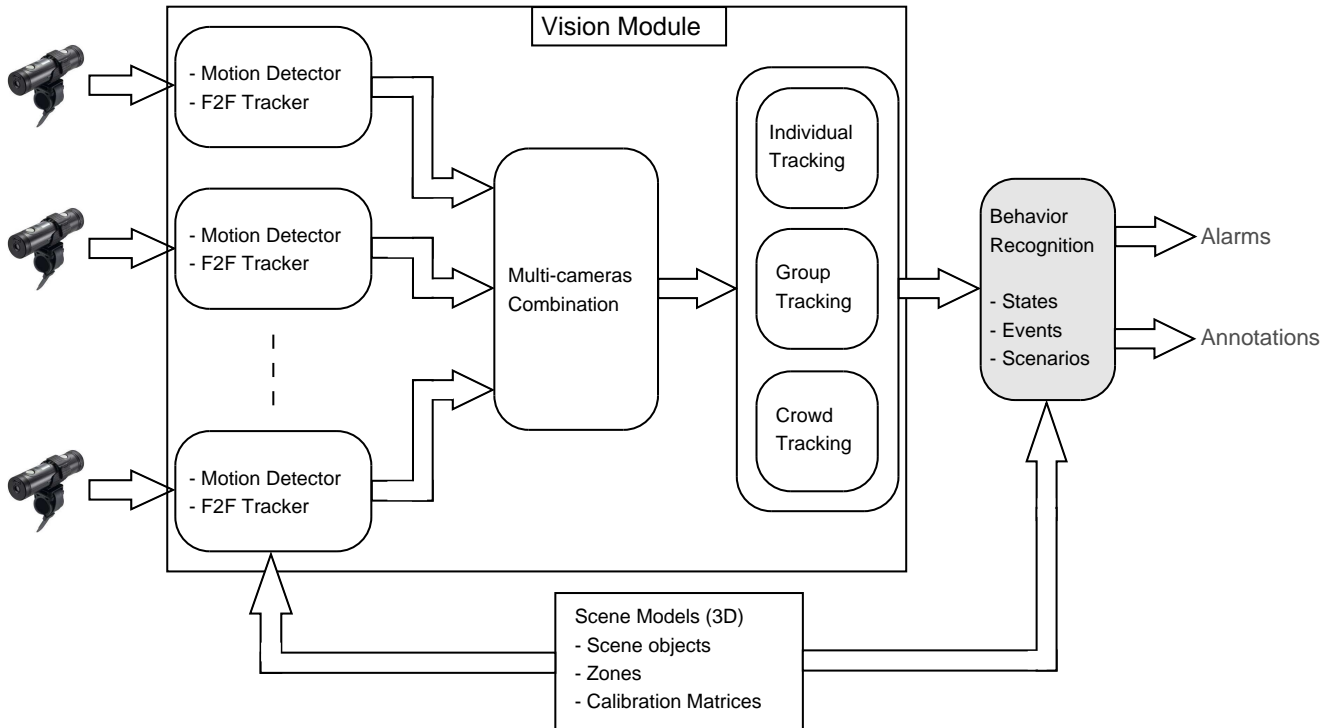


Figure 2. Video interpretation system.

entity but they are specialized depending on different configurations. For example, to compute the scenario “Fighting”, there are four methods (as shown on Figure 3). For example, one method computes the evolution of the lateral distance between people inside a group. A second one detects if someone, surrounded by people, has fallen on the floor.

4. **Output:** Contains the result of the entity computation accessible by all the other *operators*. This result corresponds to the value of the entity at the current time.

This generic framework, based on the definition of *operators*, gives two advantages: It first enables us to test a set of methods to compute an entity, independently of other entities. So we can locally modify the system (the methods to compute an entity) while keeping it globally consistent (without modifying the meaning of the entity). Second, the network of *operators* to recognize one scenario is organized as a hierarchy. The bottom of the hierarchy is composed of states and the top corresponds to the scenario to be recognized. Several intermediate levels, composed of state(s) or event(s) can be defined.

### Behavior recognition

We have defined four types of methods depending on the type of entities:

1. **Basic properties methods:** We use dedicated routines to compute properties characterizing actors such as trajectory, speed and direction. For example, we use a polygonal approximation to compute the trajectory of an individual or a group of people.

2. **State methods:** We use numerical methods which include the computation of: (a) 3D distance for states dealing with spatial relations (e.g. “an individual is close to the ticket vending machine”), (b) evolution of temporal features for states dealing with temporal relations (e.g. “the size of a group of people is constant”), (c) speed for states dealing with spatio-temporal relations (e.g. “an individual is walking”), and (d) the combination of sub states computed by other operators. The output of these numerical methods is then classified to obtain a symbolic value.

3. **Event methods:** We compare the status of states at two consecutive instants. The output of an event method is boolean: the event is either detected or not detected. For example, the event “a group of people enters a zone of interest” is detected when the state “a group of people is inside a zone of interest” changes from false to true.

4. **Scenario methods:** For simple scenarios (composed of only one state), we verify that a state has been detected during a predefined time period using a temporal filter. For sequential scenarios (composed of a sequence of states), we use finite state automaton. An automaton state corresponds to a state and a transition to an event. An automaton state also corresponds to an intermediate stage before the complete recognition of the scenario. We have used an automaton to recognize the scenarios “Blocking” and “Fraud” as described on Figure 4 and Figure 5.

For composed scenarios defining a single unit of movement composed of sub scenarios, we use Bayesian networks as proposed by (Hongeng & Nevatia, 2001) or AND/OR trees

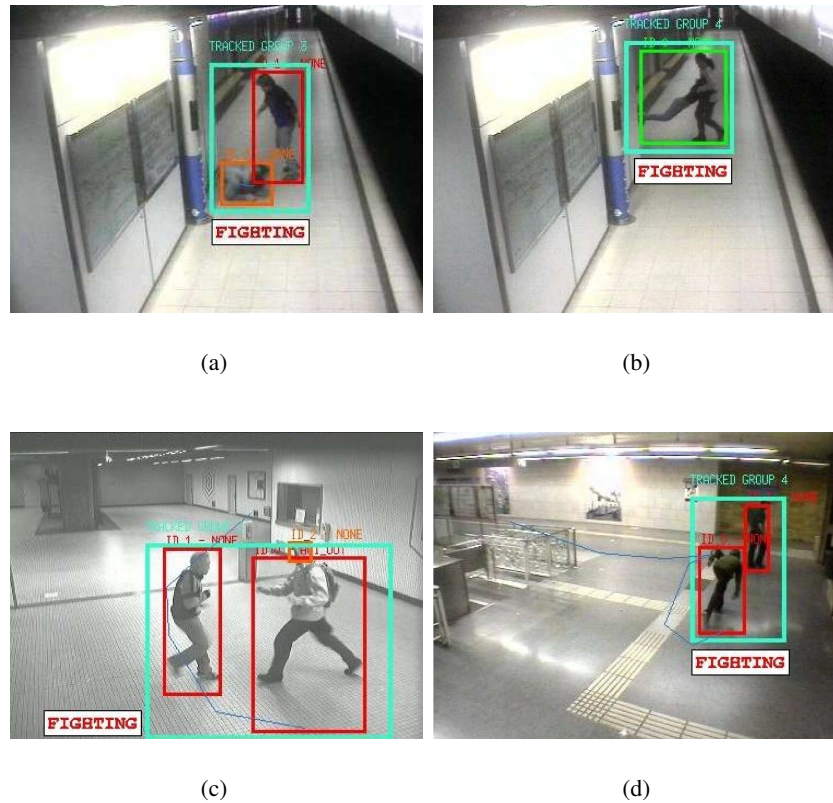


Figure 3. Four methods combined by an AND/OR tree to recognize the behavior “Fighting”. Each image illustrates a configuration where one method is more appropriate to recognize the behavior: (a) lying person on the floor surrounded by people, (b) significant variation of the group width, (c) quick separation of people inside a group, and (d) significant variation of the group trajectory.

of sub scenarios as illustrated on Figure 6. The structure of AND/OR trees, even if not continuous, is a good compromise between the usability of knowledge representation by experts and correspondence with observations on videos. A description of Bayesian networks for scenario recognition can be found in (Moenne-Loquez, Brémond, & Thonnat, 2003). We have defined one Bayesian network to recognize the “Violence” behavior composed of two sub scenarios: “Internal Violence” (e.g. erratic motion of people inside a group) and “External Violence” (e.g. quick evolution of the trajectory of the group). The structures of Bayesian networks are statistically learned by an off-line process, allowing adaptability for different kind of behaviors, but lacking in usage of knowledge from an expert. In contrast, AND/OR trees can represent more precise knowledge from experts, but not necessarily in correspondence with the observed videos. Both methods are time demanding, either to collect representative videos or tuning the parameters corresponding to the expert knowledge. Also, both need a learning stage (statistical or manual) to adjust the parameters of the network using ground truth (videos annotated by operators). Bayesian networks are optimal given ground truth but AND/OR trees are easier to tune and to adapt to new scenes.

For scenarios with multiple actors involved in complex temporal relationships, we use a network of temporal variables representing sub scenarios and we backtrack temporal constraints among the already recognized sub scenarios as proposed by (Vu et al., 2003).

For users to be able of defining the behaviors they want to recognize, an event description language is used as a formalism for describing the events characterizing these behaviors (Vu et al., 2003). The purpose of this language is to give a formal but also intuitive, easily understandable, and simple tool for describing events. All these features can be achieved by defining events in a hierarchical way and reusing definitions of simple events in more complex ones. A definition of an event consists of:

1. **Event name.**
2. **List of physical objects involved in the event:** A *physical object* can be a mobile object or a static one. Typical examples are humans, vehicles, zones or equipments.
3. **List of components representing sub events which describe simpler activities.**
4. **List of constraints expressing additional conditions:** The *constraints* can be *spatial* or *temporal* in dependence on their meaning. In both cases we can have *symbolic* or *nu-*

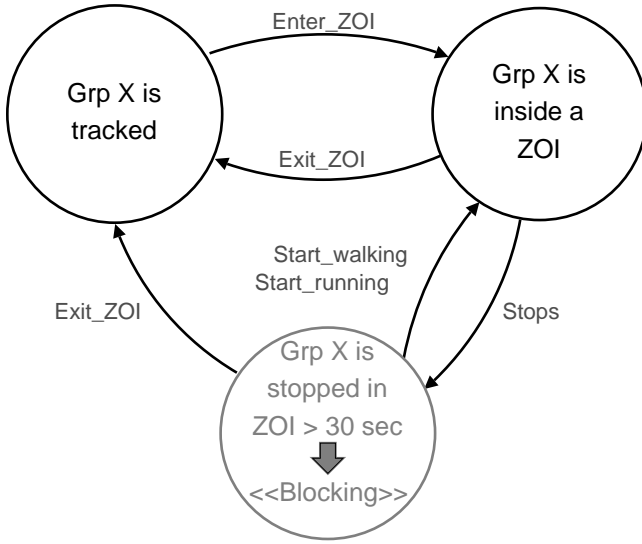


Figure 4. To check whether a group of people is blocking a zone of interest (ZOI), we have defined an automaton with three states: (a) a group is tracked, (b) the group is inside the ZOI, and (c) the group has stopped inside the ZOI for at least 30 seconds.

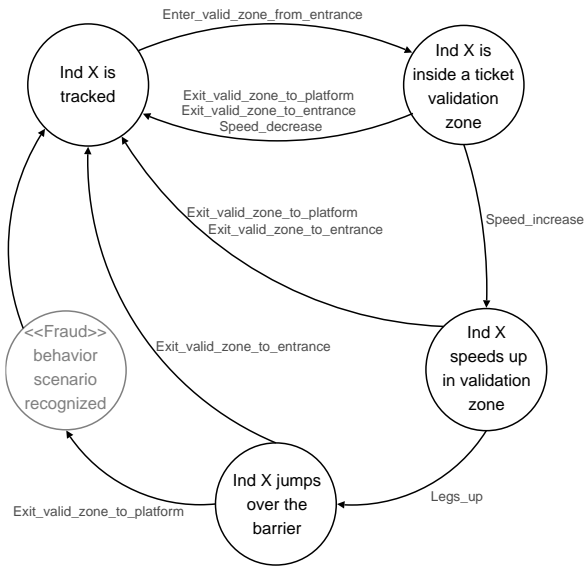


Figure 5. To check whether an individual is jumping over the barrier without validating his ticket, we have defined an automaton with five states: (a) an individual is tracked, (b) the individual is at the beginning of the validation zone, (c) the individual has a high speed, (d) the individual is over the barrier with legs up, and (e) the individual is at the end of the validation zone.

```

event ( vandalism,
  physical_objects( p: Person, eq: Equipment),
  components( (state e1: p far_from eq ),
    (state e2: p stays_at eq ),
    (event e3: p moves_away_from eq ),
    (event e4: p moves_close_to eq ),
    (state e5: p stays_at eq )
  ),
  constraints( ( (e1 before e2) (e2 before e3)
    (e3 before e4) (e4 before e5)
  )
)
)

```

Figure 7. Definition of the behavior “vandalism”. A sequence of five events which represent relative positions between a person  $p$  and an equipment  $eq$ , must be detected in order to recognize this behavior.

meric form. For example, a *spatial symbolic constraint* is “object inside zone”, while a *spatial numeric constraint* can be defined as follows:

$$\text{distance}(\text{object}_1, \text{object}_2) \leq \text{threshold}$$

In the case of a *temporal constraint*, we can also have a *numeric form* like:

$$\text{duration}(\text{event}) \leq 20[\text{secs}]$$

or a symbolic form:

**event<sub>1</sub> before event<sub>2</sub>**

On figure 7 is depicted an example of a complex scenario, “vandalism”: a person  $p$  tries to *break up* an equipment  $eq$ , using the formalism of (Vu et al., 2003). This scenario will be recognized if a sequence of five events described on figure 8 has been detected.

### Behavior recognition results

The platform has been tested in different situations and validated in the metro monitoring application. The behavior recognition module is running on a PC Linux and is processing four tracking outputs corresponding to four cameras with a rate of five images per second. We have tested the whole video interpretation system (including motion detection, tracking and behavior recognition) on videos coming from ten cameras of Barcelona and Brussels metros. We correctly recognized the scenario “Fraud” 6/6 (six times out of six) (Figure 9(a)), the scenario “Vandalism” 4/4 (Figure 9(b)), the scenario “Fighting” 20/24 (Figure 3), the scenario “Blocking” 13/13 (Figure 9(c)) and the scenario “Overcrowding” 2/2 (Figure 9(d)). We also tested the system over long sequences (10 hours) to check the robustness over false alarms. For each behavior, the rate of false alarm is: two for “Fraud”, zero for “Vandalism”, four for “Fighting”, one for “Blocking” and zero for “Overcrowding”.

Moreover, in the framework of the European project ADVISOR, the video interpretation system has been ported on Windows and installed at Barcelona metro in March 2003 to

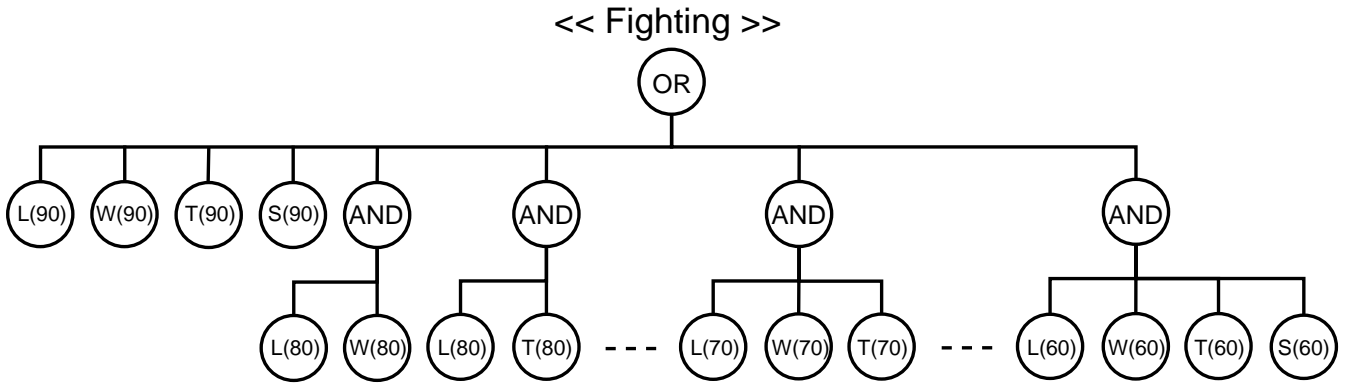


Figure 6. To recognize whether a group of people is fighting, we have defined an AND/OR tree composed of four basic scenarios: (L) lying person on the floor surrounded by people, (W) significant variation of the group width, (S) quick separation of people inside the group, and (T) significant variation of the group trajectory. Given these four basic scenarios we were able to build an OR node with all combinations (corresponding to 15 sub scenarios) of the basic scenarios. These combinations correspond to AND nodes with one up to four basic scenarios. The more basic scenarios there are in AND nodes, the less strict is the recognition threshold of each basic scenario. For example, when there is only one basic scenario (e.g. L(90)), the threshold is 90 and when there are four basic scenarios, the threshold is 60. To parameterize these thresholds, we have performed a learning stage consisting in a statistical analysis of the recognition of each basic scenario.

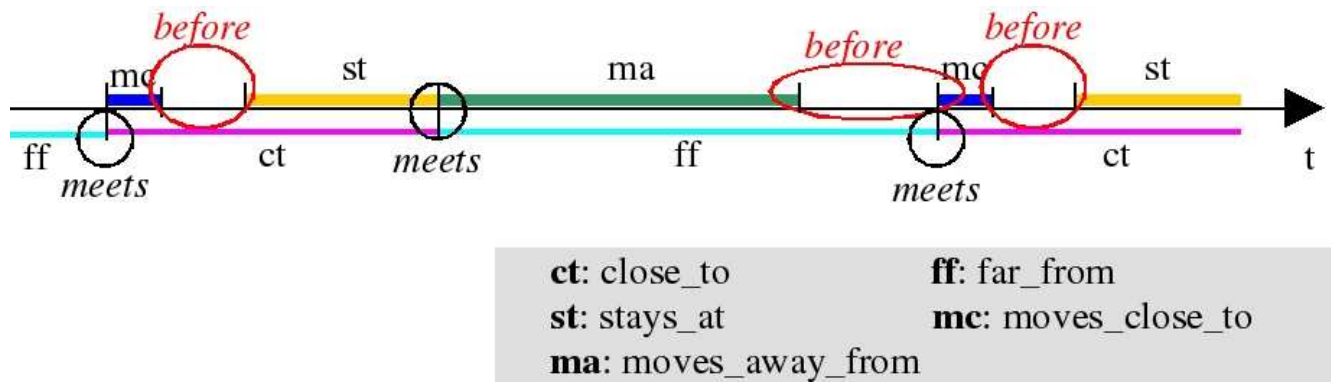


Figure 8. Temporal constraints for the states and events constituting a scenario “vandalism”.

Table 1

Results of the technical validation of the metro monitoring system. For each scenario, we report in particular the percentage of recognized instances of this scenario (fourth column) and the accuracy in time of the recognition (the percentage of the duration of the shown behavior *covered* by the generation of the corresponding alert by the system). This value is an average over all the scenario instances (fifth column).

Scenario Name	Number of Behaviors	Number of Recognized Instances	% of Recognized Instances	Accuracy	Number of False Alerts
Fighting	21	20	95 %	61 %	0
Blocking	9	7	78 %	60 %	1
Vandalism	2	2	100 %	71 %	0
Jumping o.t.b.	42	37	88 %	100 %	0
Overcrowding	7	7	100 %	80 %	0
<b>TOTAL</b>	<b>81</b>	<b>73</b>	<b>90 %</b>	<b>85 %</b>	<b>1</b>



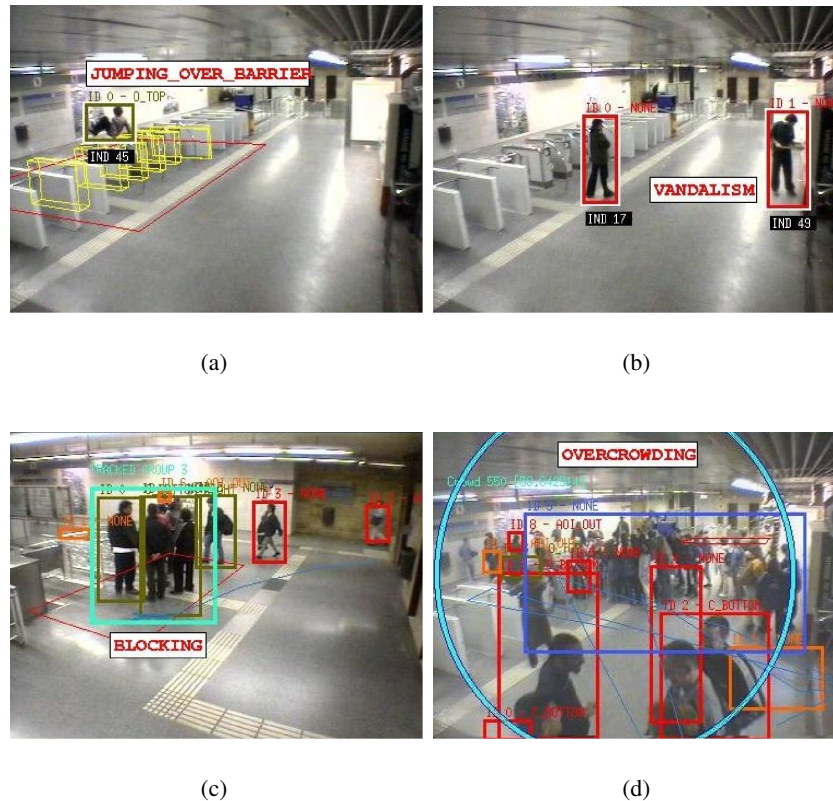


Figure 9. Four behaviors selected by end users and recognized by the video interpretation system: (a) “Fraud” recognized by an automaton, (b) “Vandalism” recognized by a temporal constraint network, (c) “Blocking” recognized by an automaton, and (d) “Overcrowding” recognized by an AND/OR tree.

be evaluated and validated. This evaluation has been done by Barcelona and Brussels video-surveillance metro operators during one week at the Sagrada Familia metro station. Together with this evaluation, a demonstration has been performed to various guests, including the European Commission, project Reviewers and representative of Brussels and Barcelona Metro to validate the system. The evaluation and the demonstration were conducted using both live and recorded videos: four channels in parallel composed of three recorded sequences and one live input stream from the main hall of the station. The recorded sequences enabled to test the system with rare scenarios of interest (e.g. fighting, jumping over the barrier, vandalism) whereas the live camera allowed to evaluate the system against scenarios which often happen (e.g. overcrowding) and which occurred during the demonstration and also to evaluate the system against false alarms. In total, out of 21 fighting incidents in all the recorded sequences, 20 alarms were correctly generated, giving a very good detection rate of 95%. Out of nine blocking incidents, seven alarms were generated, giving a detection rate of 78%. Out of 42 instances of jumping over the barrier, including repeated incidents, the behavior was detected 37 times, giving a success rate of 88%. The two sequences of vandalism were always detected over the six instances of vandalism, giving

a perfect detection rate of 100%. Finally, the overcrowding incidents were also consistently detected in the live camera, with some 28 separate events being well detected. Results are summarized in Table 1.

In conclusion, the ADVISOR demonstration has been evaluated positively by end-users and European Committee. The algorithm responded successfully to the input data, with high detection rates, less than 5% of false alarms and with all the reports being above 70% accurate.

## DISCUSSION

The evaluation of the metro application has validated the ability of VSIP for the recognition of human behaviors. End-users of the application evaluated it positively, because of its high detection rate on different scenarios.

We have been also working closely with end-users on other application domains. For example, we have built with VSIP six systems which have been validated by end-users:

1. The activity monitoring system in metro stations has been validated by metro operators from Barcelona and Brussels.

2. A system for detecting abnormal behaviors inside moving trains (Figure 10(a)), able to handle situations in which people are partially occluded by the train equipment (like

seats) has been validated over three scenarios (graffiti, theft, begging).

3. A bank agency monitoring system has been installed and validated in four bank agencies around Paris (Georis, Mazière, Brémond, & Thonnat, 2004) (Figure 10(b) and Figure 10(c)).

4. A lock chamber access control system for buildings security has been validated on more than 140 recorded videos and on a live prototype (Figure 10(d)).

5. An apron monitoring system has been developed for an airport<sup>2</sup> where vehicles of various types are evolving in a cluttered scene (Figure 10(e) and Figure 10(f)). The dedicated system has been validated on five servicing scenarios (GPU vehicle arrival, parking and refueling aircraft, loading/unloading container, towing aircraft).

6. A metro access control system has been tested by end-users on more than 200 videos and on a live prototype (Bui Ngoc, Brémond, Thonnat, & Faure, 2005) (Figure 10(g)).

Some of these applications are illustrated on Figure 10. They present several characteristics which make them interesting for research purposes:

1. The observed scenes vary from large open spaces (like metro halls) to small and closed spaces (corridors and lock chambers).

2. Cameras can have both non overlapping (like in metro stations and lock chambers systems) and overlapping fields of view (metro stations and bank agencies).

3. Humans can interact with the equipment (like ticket vending machines or access control barriers, bank safes and lock chambers doors) either in simple ways (*open/close*) or in more complex ones (such as the interaction occurring during *vandalism against equipment* or *jumping over the barrier* scenarios).

We are currently building with VSIP various other applications. For instance, a system concerning traffic monitoring on highway has been built in few weeks to show the adaptability of the platform (see Figure 10(h)). Other applications are envisaged such as home-care (monitoring of elderly people at home), multimedia (e.g. intelligent camera for video conferencing) and animal behavior recognition (e.g. insect parasitoids attacking their hosts). VSIP platform is currently extended to learn behavior models using unsupervised learning techniques to be applied on parking lot monitoring (Toshev et al., 2006) (see Figure 10(i) and Figure 10(j)).

VSIP has shown its ability to automatically recognize and analyze human behaviors. However, some limitations remain. Video understanding systems are often difficult to configure and install. To have an efficient system handling the variety of the real world, extended validation is needed. Automatic capability to adapt to dynamic environments should be added to the platform, which is a new topic of research. Nevertheless, the diversity of applications where VSIP has been used shows that this platform is suitable to fulfill many types of requirements.

## CONCLUSION

We have presented a video understanding platform to automatically recognize human behaviors involving individuals, groups of people, crowds and vehicles, by detecting visual invariants. A visual invariant is a visual property or clue which characterizes a behavior.

Different methods have been defined to recognize specific types of behaviors under different scene configurations. To also experiment various software solutions, all these methods have been integrated in a coherent framework enabling to modify locally and easily a given method. VSIP platform has been fully evaluated on several applications. For instance, the system has been tested off-line and has been evaluated, demonstrated and successfully validated in live condition during one week at the Barcelona metro, in March 2003.

This approach can also be applied on biological domain. For instance, in 2005 we are developing a system based on VSIP platform which detects the behaviors of a trichogramma interacting with butterfly eggs. This application corresponds to a behavioral study for understanding how a trichogramma introduces its eggs on butterfly eggs, contributing to plague control on agriculture.

Hence, we believe that VSIP shows a great potential as a tool for recognition and analysis of human behaviors in very different configurations.

VSIP still presents some limitations when environmental conditions suddenly change or complexity of the scene increases, which makes necessary the improvement of vision modules to ensure robustness. Moreover, VSIP requires the pre-definition of events for the detection of behaviors, which can be a very hard task. To cope with this limitation, an unsupervised frequent events detection algorithm (Toshev et al., 2006) has shown encouraging preliminary results. This algorithm is capable of extracting the most significant events in a scene, without behavior pre-definition by end-user.

## References

- AVITRACK European Research Project, <http://www.avitrack.net>
- Bobick, A. F., & Wilson, A. D. (1997). A state-based approach to the representation and recognition of gesture. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 19(12), 1325-1337.
- Borg, M., Thirde, D., Ferryman, J., Fusier, F., Brémond, F., & Thonnat, M. (2005, September). Video event recognition for aircraft activity monitoring. In *Proceedings of the 8<sup>th</sup> International IEEE Conference on Intelligent Transportation Systems (ITSC05)*. Vienna, Austria: IEEE Computer Society.
- Brand, M., & Kettner, V. (2000). Discovery and segmentation of activities in video. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 22(8), 844-851.
- Brémond, F., Maillot, N., Thonnat, M., & Vu, T. V. (2004, May). *RR5189 - Ontologies for video events* (Tech. Rep.). Sophia Antipolis, France: Orion Team, Institut National de Recherche en Informatique et Automatique (INRIA).

<sup>2</sup> In the framework of AVITRACK European Project, see (AVITRACK) and (Borg et al., 2005)

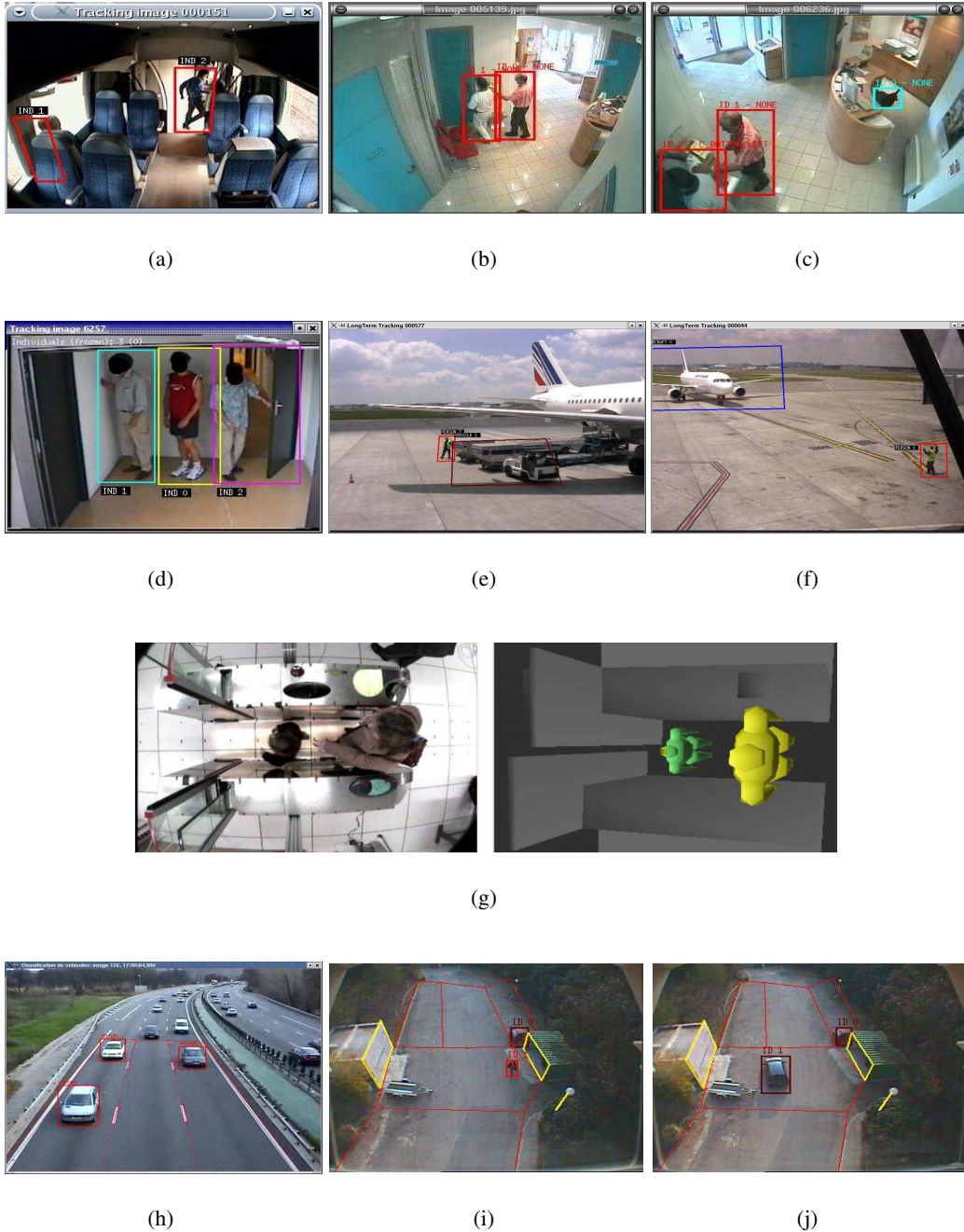


Figure 10. Six other systems derived from VSIP to handle different applications. Image (a) illustrates a system for unruly behaviors detection inside trains. Images (b) and (c), taken with two synchronized cameras with overlapping field of view working in a cooperative way, depict a bank agency monitoring system detecting an abnormal *bank attack* scenario. Image (d) illustrates a single-camera system for a lock chamber access control application for building entrances. Images (e) and (f) show a system for apron monitoring on airports; this system combines a total of eight surveillance cameras with overlapped fields of view. Image (g) depicts a multi-sensor system for metro access control; this system uses the information obtained from a static top camera and a set of lateral visual sensors, to perform shape recognition of people. The right image corresponds to the understanding of the scene projected in a 3D virtual representation. Image (h) illustrates a highway traffic monitoring application. Finally, images (i) and (j) depict a parking lot monitoring system for learning behavior patterns using an unsupervised technique.

- Bui Ngoc, H.-B., Brémond, F., Thonnat, M., & Faure, J.-C. (2005, September). Shape recognition based on a video and multi-sensor system. In *Proceedings of the IEEE International Conference on Advanced Video and Signal-Based Surveillance (AVSS 2005)* (p. 230-235). Los Alamitos, CA: IEEE Computer Society Press.
- Chleq, N., & Thonnat, M. (1996, September). Real-time image sequence interpretation for video-surveillance applications. In *Proceedings of the IEEE International Conference on Image Processing (ICIP96)* (Vol. 2, p. 801-804). Lausanne, Switzerland: IEEE Computer Society.
- Cupillard, F., Brémond, F., & Thonnat, M. (2004, 21-23 March). Video understanding for metro surveillance. In *Proceedings of the IEEE International Conference on Networking, Sensing and Control (ICNSC04)* (Vol. 1, p. 186-191). Taipei, Taiwan: IEEE Computer Society.
- Galata, A., Cohn, A., Magee, D., & Hogg, D. (2002). Modeling interaction using learnt qualitative spatio-temporal relations and variable length markov models. In *Proceedings of the 15<sup>th</sup> European Conference on Artificial Intelligence (ECAI02)* (p. 741-745). Lyon, France: IOS Press.
- Georis, B., Mazière, M., Brémond, F., & Thonnat, M. (2004, February). A video interpretation platform applied to bank agency monitoring. In *Proceedings of the International Conference on Intelligent Distributed Surveillance Systems (IDSS04)* (p. 46-50). London, Great Britain: The Institution of Electrical Engineers.
- Gerber, R., Nagel, H., & Schreiber, H. (2002, 21-26 July). Deriving textual descriptions of road traffic queues from video sequences. In *Proceedings of the 15<sup>th</sup> European Conference on Artificial Intelligence (ECAI02)* (p. 736-740). Lyon, France: IOS Press.
- Ghallab, M. (1996, 5-8 November). On chronicles: Representation, on-line recognition and learning. In *Proceedings of the 5<sup>th</sup> International Conference on Principles of Knowledge Representation and Reasoning (KR96)* (p. 597-606). Cambridge, Massachusetts, USA: Morgan Kaufmann.
- Hongeng, S., Brémond, F., & Nevatia, R. (2000). Representation and optimal recognition of human activities. In *IEEE Proceedings of the International Conference on Computer Vision and Pattern Recognition (CVPR00)*. South Carolina, USA: IEEE Computer Society.
- Hongeng, S., & Nevatia, R. (2001, May). Multi-agent event recognition. In *Proceedings of the 8<sup>th</sup> International Conference on Computer Vision (ICCV2001)*. Vancouver, B.C., Canada: IEEE Computer Society.
- Howell, A., & Buxton, H. (2002, October). Active vision techniques for visually mediated interaction. *Image & Vision Computing*, 20(12), 861-871.
- Hu, W., Tan, T., Wang, L., & Maybank, S. (2004). A survey on visual surveillance of object motion and behaviors. *IEEE Transactions on Systems, Man, & Cybernetics - Part C: Applications & Reviews*, 34(3), 334-352.
- Ivanov, Y. A., & Bobick, A. F. (2000). Recognition of visual activities and interactions by stochastic parsing. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 22(8), 852-872.
- Moenne-Loquez, N., Brémond, F., & Thonnat, M. (2003, April). Recurrent bayesian network for the recognition of human behaviors from video. In *Proceedings of the 3<sup>rd</sup> International Conference on Computer Vision Systems (ICVS03)*. Graz, Austria: Springer.
- Owens, J., & Hunter, A. (2000, July). Application of the self-organizing map to trajectory classification. In *Proceedings of the IEEE International Workshop on Visual Surveillance (VS2000)*. Dublin, Ireland: IEEE Computer Society.
- Pinhanez, C., & Bobick, A. (1998, June). Human action detection using pnf propagation of temporal constraints. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR98)* (p. 898-904). Santa Barbara, CA, USA: IEEE Computer Society.
- Rota, N., & Thonnat, M. (2000, August). Activity recognition from video sequences using declarative models. In *Proceedings of the 14<sup>th</sup> European Conference on Artificial Intelligence (ECAI00)*. Berlin, Germany: IOS Press.
- Toshev, A., Brémond, F., & Thonnat, M. (2006, January). Unsupervised learning of scenario models in the context of video surveillance. In *Proceedings of the IEEE International Conference on Computer Vision Systems (ICVS06)* (p. 10). New York, NY, USA: IEEE Computer Society.
- Vu, T.-V., Brémond, F., & Thonnat, M. (2003, August). Automatic video interpretation: a novel algorithm for temporal scenario recognition. In *Proceedings of the 18<sup>th</sup> International Joint Conference on Artificial Intelligence (IJCAI03)*. Acapulco, Mexico: Morgan Kaufmann.
- Xie, L., Chang, S.-F., Divakaran, A., & Sun, H. (2003). Unsupervised mining of statistical temporal structures in video. In A. Rosenfeld, D. Doermann, & D. Dementhon (Eds.), *Video mining* (p. 279-307). Norwell, MA, USA: Kluwer Academic.