

A Framework for Activity Recognition and Detection of Unusual Activities

Dhruv Mahajan Nipun Kwatra Sumit Jain Prem Kalra
Subhashis Banerjee

Department of Computer Science and Engineering
Indian Institute of Technology
New Delhi 110016
email: {pkalra,suban}@cse.iitd.ernet.in

Abstract

In this paper we present a simple framework for activity recognition based on a model of multi-layered finite state machines, built on top of a low level image processing module for spatio-temporal detections and limited object identification. The finite state machine network learns, in an unsupervised mode, usual patterns of activities in a scene over long periods of time. Then, in the recognition phase, usual activities are accepted as normal and deviant activity patterns are flagged as abnormal. Results, on real image sequences, demonstrate the robustness of the framework.

1. Introduction

Analyzing complex spatio-temporal changes in a dynamic scene in order to recognize logical sequences of usual activity patterns and detect unusual deviations has become an important problem in computer vision, especially in the context of visual surveillance. In this paper we present a simple framework for activity recognition based on a model of multi-layered finite state machines, built on top of a low level image processing module for spatio-temporal detections and limited object identification. We fix the modes of interaction in the lower physical layer of the network architecture depending on the context of the problem, and automatically find clusters of high probability sequences of transitions to learn usual behavioural patterns in an unsupervised manner. Low probability sequences which are not recognized by the FSM network are diagnosed as unusual. We argue that such a model is more suitable for detection and recognition of activity patterns in surveillance situations than some of the more complex models suggested in the literature - based on Hidden Markov Models [16, 12, 17, 14, 15], Bayesian networks [2, 13, 8, 9] and stochastic context free grammars [11].

In contrast with most of the techniques found in the literature which are either hard coded to recognize only specific activities or model activity patterns using supervised learn-

ing, we consider two different scenarios -

1. unsupervised learning of usual activity patterns and detection of unusual activities. Any activity which is not recognized as normal is flagged as deviant.
2. explicit programming of the FSM's for recognition of complex activities or training using supervised learning.

Examples of the first of the above may include learning usual activity patterns of movements of people in an airport terminal or of people and cars in a parking lot, and detecting deviant patterns like walking in a wrong direction, straying in to restricted areas, crossing of barriers, wrong parking and suspicious behaviours like taking more than usual amount of time to open the door of a car. Our framework for this case is generic - only the low level detectors in the image processing layer and the physical layer of the network structure need to be changed with context, the algorithms for activity learning and recognition require no changes whatsoever and can be trained using unsupervised learning, even for multiple types of unknown activities.

Examples of the second may include *unattended baggage detection*, wherein a person leaves a bag in the scene and walks away and the bag lies unattended for a certain duration of time; and *stealing*, where a person puts a bag down and a second person picks it up and rushes away. Our system can either be explicitly programmed to recognize activities such as these or may be trained to recognize these activities by enacting them a number of times.

The rest of the paper is organized as follows. In Section 2 we review some of the literature on activity recognition and detection and put our approach in perspective. In Section 3 we describe our low level image processing module and object detectors. In Section 4 we describe our activity analysis framework based on layered finite state machines. We present some results on examples of the type mentioned above in Section 5. Finally, in Section 6, we conclude the paper.

2. Related work

Recently, the problem of activity detection and recognition in the context of visual surveillance has received considerable attention [1]. There has been significant work that span across techniques for low level event detection [5, 20, 23, 4] and activity modeling [16, 12, 17, 14, 15, 11, 2, 13, 8, 9, 6, 21, 3, 10].

Starting from the early work of Yamato et. al. [22], HMM's have been a popular tool for activity modeling, motivated primarily by its successful use in speech recognition. A HMM is a stochastic finite state machine which models an activity pattern by learning transition probabilities among its non-observable states such that the likelihood of observation of a temporal sequence of symbols representing the activity is maximized. HMM's have been used to model simple isolated hand gestures [19], and more complex variants like coupled HMM's [17] and layered HMM's [16] have been proposed to model events such as interaction between multiple mobile objects. Some of these enhancements also implicitly model dynamic time warping. However, unlike in speech (and perhaps in gestures, to a certain extent), where the relationship of the observed output symbol to the internal semantic state is uncertain at any instant of time necessitating the use of 'hidden' states, in surveillance applications the relationship is often more direct. Consider, for example, the stealing example (see above). The internal semantic states are clearly characterized by the sequence 'person A puts a bag down', 'person B approaches the bag', 'person B picks up the bag', 'person B rushes away' - and all of these are directly observable from the image sequence. For representing such activities it is perhaps prudent to use a model where the state representations are 'transparent' with well identified physical meanings rather than 'hidden'. This would not only facilitate direct estimation of the state transition probabilities from the observables, thereby eliminating the need for solving a complex maximum likelihood estimation problem using iterative techniques, but also avoid the problems of over-parametrization and over-fitting especially when the training data is inadequate.

In [11] a parsing technique based on stochastic context-free grammars is proposed to compute the probability of a temporally consistent sequence of primitive actions recognized by a HMM model. A model that recognizes a context-free grammar is indeed more powerful than a FSM model which can capture only regular languages. For example, it allows for recognition of languages like $\{a^n b^n, n \geq 0\}$ for logical symbols a and b , which cannot be recognized by FSM's. However, in typical surveillance applications, activity patterns that need to be modeled are seldom of the type that requires bracket matching using a stack, and simple precedence requirements usually suffice in most cases, which can be effectively modeled by FSM's. In [18] it is

argued that a finite state machine cannot model temporal inter-leavings of low level events which may occur concurrently and do not have simple temporal precedence relationships in complex multi-actor activities. To deal with such situations the authors propose an multi-layered extension of FSM's called *Propagation networks*. However, we show that such situations can be effectively captured by coupled FSM's.

Another popular approach for activity recognition is though the use of Bayesian networks. In fact, in an approach closely related to ours Hongeng et. al. [8, 9] propose a multi-layered FSM model for activity recognition. However, they estimate the parameters of their model using supervised training using a Bayesian formulation. The major difficulty with the Bayesian formulation lies in obtaining the apriori estimates, especially for unusual activities which are not part of a closed set. Moreover, in their model the states of the FSM's depend on absolute time clocks causing an explosion of possibilities for state changes. Also, the method cannot be directly extended to unsupervised clustering.

Recently, there have been some interesting approaches to detection of unusual activities. In [21, 3] moving objects on a 2D plane are treated as point clusters, whose deformable shapes in time represent usual patterns of activities. Deviations are then detected from the shape changes. In [10], usual activities in a video stream are learned using unsupervised clustering. However, the application scenarios of these approaches are different from ours.

In contrast to some of the other finite state models, the states in our multi-layer FSM framework have well identified physical meanings. They represent single or multiple objects, their positions and distance vectors. There is no explicit need for fixing the size of the network apriori. Also, the state transition probabilities, at any layer, can be directly estimated from the observables. As our results clearly demonstrate, our framework can also handle complex activities involving concurrent and interleaved actions by multiple actors in a scene. Further the system is totally programmable in the sense that one can feed in an FSM to the system, described on the physical states provided by the lower level image processing layer. The richer this layer is, the more the variety of FSM's can be specified on it. Most of the existing systems do not support this programmability.

3. Image analysis modules

As mentioned above, our state descriptions are derived from physically observable attributes like object types and their positions on the ground plane. We use standard image processing techniques to detect and compute these from on-line sequences [5]. Our system works in real-time at full frame rate. We assume the ground plane to be calibrated in terms

of a XY grid whose resolution is dependent on the context. We also assume that the homography of the ground plane to be known from calibration [7]. Our low level image processing layer consists of the following modules:

1. segmentation using adaptive multi-layer background subtraction [5].
2. identification of splitting and merging of image segments.
3. object identification using principle component analysis. Currently, we detect only three types of objects - humans, cars and normal bags.
4. object tracking using Kalman filters and inter-frame object matching using volume intersection of hue-saturation histograms of detected objects.
5. detection of the position of the object on the XY grid on the ground plane. For humans we compute this by identifying the position of the feet on the ground plane as the lowest point in the segmented region along the vertical direction (whose vanishing point is known from calibration) and then applying the homography of the ground plane. For top views of cars we compute the centroid of the segmented region and apply the homography. For perspective views of cars we locate the lowest point.

Of course, segmentation is problematic to say the least, and the low level image processing is clearly the weakest link of our method. In the present state of our implementation dense situations, severe occlusions and unfavourable lighting result in failures more often than not. Consequently, our experiments are mainly restricted to sparse situations and good illumination conditions.

4. Framework of multi-layered finite state machines

Our architecture of multi-layered FSM's for learning and recognition of complex actions is depicted in Figure 1.

The low-level image processing methods described in the previous section can provide only a physical description of the dynamic scene. The semantic interpretation of the physical events is carried out in the FSM's in the *physical layer* of the architecture. The interpretation at this level is in terms of the number of objects of different types present in the scene, their split and merge history, their geometric positions, and changes in these positions with time. High probability sequences of single object motions and multi-object interactions at the physical layer represent logical interpretations of the physical actions. Examples of such logical interpretations may be a person walking or running in

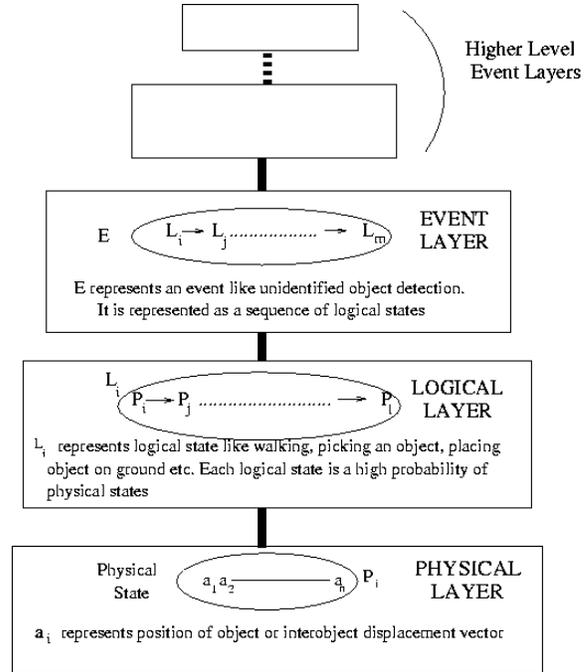


Figure 1: Multi-layered FSM model

a certain direction, a person walking toward a car, a person putting a bag down, two people moving toward each other etc. These logical actions in time represent the symbols and states in the higher *logical layer*. The types of interactions between objects that are to be considered are specified a priori depending on the context. The logical states are determined automatically as high probability sequences of physical transitions.

In many situations, where interactions between multiple objects are unimportant, the scene interpretation can be completed at the logical layer level itself. For example, consider the monitoring of proper/improper parking of individual cars, or straying of individuals in to restricted areas. Such logical events can be flagged without even considering inter-object interactions. However, for more complex actions one requires to analyze high probability sequences of such primitive logical actions. Examples include a) a person walking toward a car, opening the door of the car and merging with the car and driving away - a sequence of three high probability events in the logical layer involving a person and a car, or b) stealing as described before involving two people and a bag. The high probability sequences in the logical layers are extracted as events which are states in the *event layer*. Even more complex events which are characterized by high probability sequences of simpler events are detected at the event layer.

Note that when trained in the unsupervised mode, the states at the higher layers are just high probability sequences

at the lower layers detected automatically by clustering. While these states have clear semantics, they are still nameless. Recognition of usual behavioural patterns and detection of deviant ones can still be carried out regardless. In what follows we describe the details of each layer.

4.1. The physical layer

The structure of the physical layer is specific to the context. Let O_p be the p^{th} object detected in the scene which may be one of K possible types. For each such O_p we dynamically create a FSM at the lowest physical layer. The states of the FSM's are the geometric positions XY in the rectangular grid partitioning the ground plane. The edges in these FSM's represent transitions to neighbouring grid positions. We fix a priori what type of inter-object interactions are important from the context. If objects O_p and O_q interact, we also create a FSM corresponding to $O_p O_q$, the states of which represent all possible distance vectors from positions of O_p to positions of O_q measured in the grid coordinates, and their merging and splitting history. For example, in a parking lot scenario, we create a FSM for each $(car_p, person_q)$ pair of detected cars and people, in addition to individual FSM's for each car_p and $person_q$ (assuming that we are not interested in modeling possible interactions between two cars or two people, like accidents or collisions). Each state of a $(car_p, person_q)$ FSM is of the type $(\Delta X, \Delta Y, status)$ where $status$ is one of $\{split, merge, separate\}$. Each state of a car_p or $person_q$ FSM is of the type (X, Y) .

For each FSM corresponding to an instance of an object or inter-object interactions, we also maintain the following data structures. For M possible states, $transMx$ and $timeMx$ are $M \times M$ matrices and $time_states$ is an array of size M . $transMx[i][j]$ maintains the count of the number of transitions of the object from state i to state j . $timeMx[i][j]$ maintains the mean and the variance of the time taken for the transition. $time_state[i]$ maintains the mean and the variance of the time spent by the object in state i .

During the training phase, we aggregate all the statistics obtained from individual examples in to generic class specific FSM's. These generic FSM's accumulate the $transMx$, $timeMx$ and $time_state$ data from individual FSM's, which can be used to assign a weight P_{ij} to each edge (i, j) , where P_{ij} is the probability that given the system is in state i it makes a transition to state j . Clearly,
$$P_{ij} = \frac{transMx[i][j]}{\sum_j transMx[i][j]}.$$

Statistics on time taken for transition ($timeMx$) between states and time spent ($time_state$) in different states by the system can also be calculated in these generic FSM's. For example, in the parking lot scenario, all state transition statistics from individual FSM's like car_p , $person_q$ and

$$\begin{aligned} S &= \{i, j\} \\ P &= P_{ij} \\ seq &= \phi \\ forward &= false \end{aligned}$$

```

find_sequence(S, P, forward) =
  i = first_state(S)
  j = last_state(S)
  in_i = incoming_edges(i)
  out_j = outgoing_edges(j)
  for each e in out_j
    if Pje > Tedge ∧ P * Pje > Tseq then {
      find_sequence(S ∪ {e}, P * Pje, true)
    }
  for each e in in_i
    if Pei > Tedge ∧ P * Pei > Tseq then {
      if forward == false then
        find_sequence({e} ∪ S, P * Pei, false)
    }
  if no edge qualifies above threshold conditions then
    seq = seq ∪ S

```

Figure 2: Algorithm for finding *HPL* sequences

$(car_p, person_q)$ are aggregated in to three generic FSM's called cars, people and (car, people) respectively.

4.2. The logical layer

As we have mentioned above, we consider the high probability sequences of the physical layer FSM's to represent logical events. We want to extract all maximal sequences $S_n = (i, j, \dots, n)$ such that $P(S_n) = P_{ij}P_{jk} \dots P_{n} > T_{seq}$ and $P_{lm} > T_{edge}$ for every consecutive pairs of states (l, m) in S_n , where T_{seq} and T_{edge} are the probability thresholds for a sequence and an edge respectively. We call such sequences *HPL* sequences (High Probability Logical sequences). Starting from any edge (i, j) , the algorithm in Figure 2 finds all sequences S_n with (i, j) as an edge by examining sub-sequences in both forward and backward directions. We ignore the backward branches of a sequence generated by all forward expansions in order to avoid reaching the same sequence from more than one path, as shown in Figure 3.

Every *HPL* sequence in the physical layer represents a state in the logical layer. The system is in a logical state i if the system has most recently gone through a sequence of physical states corresponding to the logical state i . The state transitions at the logical layer happens when the physical state shifts from one *HPL* sequence to another.

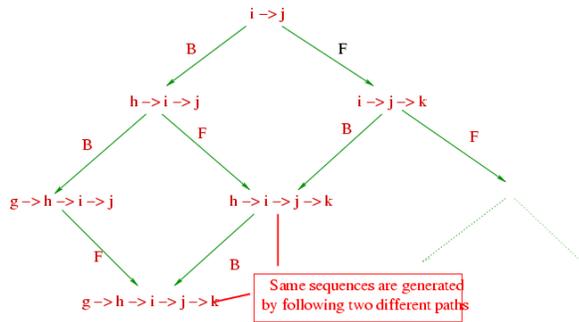


Figure 3: Example of multiple paths for the same sequence

4.3. Event layer

For determining complex events which are high probability transitions in the logical layer, we employ an identical algorithm at the logical layer. In our current experiments, we have restricted our framework to three levels of hierarchy.

4.4. Aggregation of training data

During the training phase, the data is processed on-line to do the following tasks - i) build generic FSM's in the physical states as described in Section 4.1 and ii) building of higher level states: with each transition in the physical state the HPL algorithm is run to see if any new logical states are found. Any such logical state is pushed up and stored for detection of unusual events as described later. The same process is repeated for each higher level layer.

4.5. Detection of unusual events

We maintain a counter *failcount* with each *HPL* sequence, which tracks the number of times the FSM corresponding to a *HPL* sequence fails to make a valid transition. A global counter *totalcount* maintains a count of the time from which any logical transition has occurred. If *totalcount* exceeds a certain threshold, then each *failcount* is checked. If all the *failcount*'s are over a threshold, then the system is in an undefined logical state and an unusual event is flagged. If any transition in a *HPL* sequence is faster or slower than that predicted by the learning phase statistics, then an *unusual time* alarm is flagged. This, for example, can indicate running when the usual behavioural pattern is to walk.

4.6. Explicit programming of the FSM's

Finally, an FSM network described on the physical states provided by the lower level image processing layer can just be fed into the system for detecting unusual events in a supervised manner. In Figure 10 we give one such example of the FSM at the event layer corresponding to the sequence of logical events leading to the interpretation of stealing.



Figure 4: The walking example, recognition of a usual event.

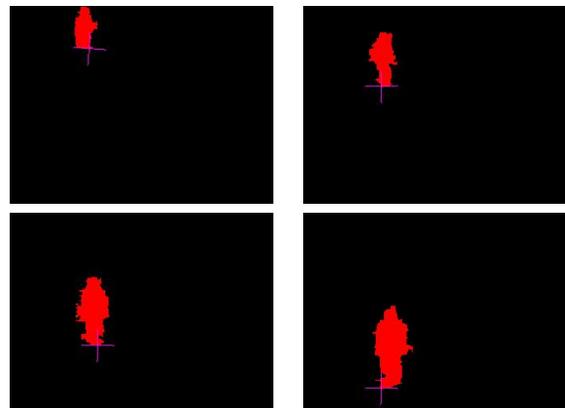


Figure 5: The walking example, back-ground subtraction and tracking.

5. Results

In what follows we present some results on activity recognition. Please visit http://www.cse.iitd.ernet.in/vglab/research/research_group2/1/activmon/index.shtml to see the videos.

In Figure 4 we show results of recognition of usual walking patterns. The system was trained (unsupervised) with sequences in which people walk in the corridor in either directions. The green box (upper right corner) indicates that some logical sequence has been completed. In Figure 5 we show the corresponding back-ground subtracted frames and the extracted ground positions. In Figure 6 we show a sequence in which the person tries to cross a barrier and jump in to thin air. The red signal (upper right corner) indicates that unusual sequence transitions have happened.

In Figure 7 we show results in a parking lot scenario. The system was trained with sequences in which different cars move in to the parking slots. The slot in front of the gate (fourth from the right) is actually marked as 'No Park-

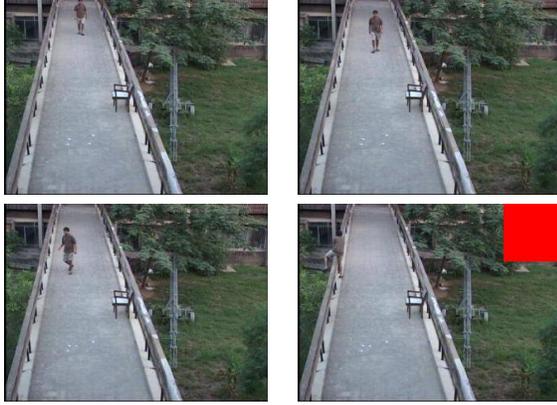


Figure 6: The walking example, detection of an unusual event.

ing’. The white car is recognized to have parked correctly, indicated by the small green square on the car. The parking of the red car in to the illegal area is flagged as an unusual event, indicated by a red square on the car. The training was unsupervised.

Both of the above examples are based on learning of transitions of absolute positions of individual objects. In the example of Figure 8 we show results of an FSM representing interaction of two objects - a person and a car. The states of the FSM represent distance vectors between a person and a car, and events such as splitting and merging. The system is trained with sequences in which people walk toward cars, open the door and get in. The figure on the left shows usual recognition - the green box indicates that some logical sequence has been completed. The figure on the right shows a frame of the sequence in which the time spent in the state is more than usual - the blue signal indicates a time alarm, i.e., system is in a particular state for more than the expected time. In Figure 9, we show a similar example where the walking people interact with several cars. Usual events recognized are ‘walking past a standing car’, ‘approaching a car’, ‘opening the door and getting in’. These are marked with a green square on the person. Unusual events detected are ‘more than usual amount of time spent’ - marked with a blue square on the person.

Our final result, in Figure 10, corresponds to explicit programming of stealing. We show one frame of a stealing sequence, and the corresponding state as detected during a run of the FSM.

6. Conclusion

We have presented a new framework for unsupervised learning of usual activity patterns and detection of unusual activities based on a network of finite state machines. The FSM framework is simple yet powerful and can reliably

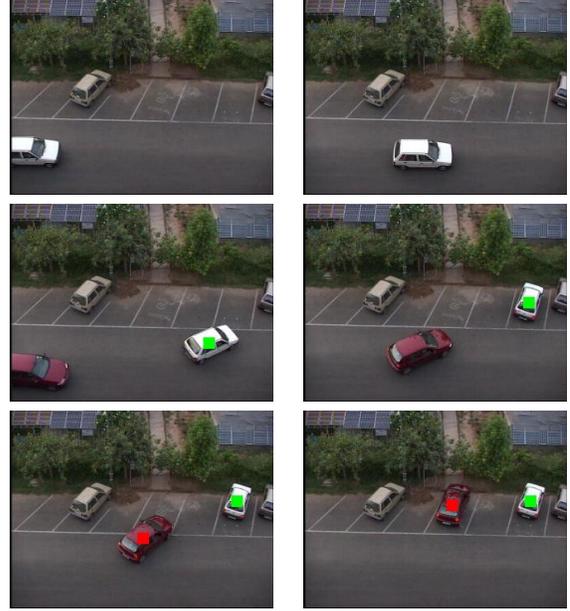


Figure 7: Frames from a car-parking example.



Figure 8: An example of detection of possible car theft.

capture complex and concurrent inter-object interactions. Our preliminary results demonstrate the potential of the approach. In our present implementation, we do not handle the cases of complex image processing and tracking under occlusion. In future we plan to include such cases and test the framework rigorously under dense situations with many activities happening simultaneously. Presently our hypothesis about the objects in the frame is based on a single frame which may not be very reliable due to image processing errors. Error correction in higher level layers can be added to make the system more robust. In future we plan to incorporate such error corrections in our system.

References

- [1] J. K. Aggarwal and Q. Cai. Human Motion Analysis: A Review. *IEEE Trans. on Pattern Anal. and Machine Intell.*, 21(11):123 – 154, 1999.
- [2] H. Buxton and S. Gong. Advanced visual surveillance using bayesian networks. In *Proc. Intl. Conf. Computer Vision*, pages 111 – 123, 1995.

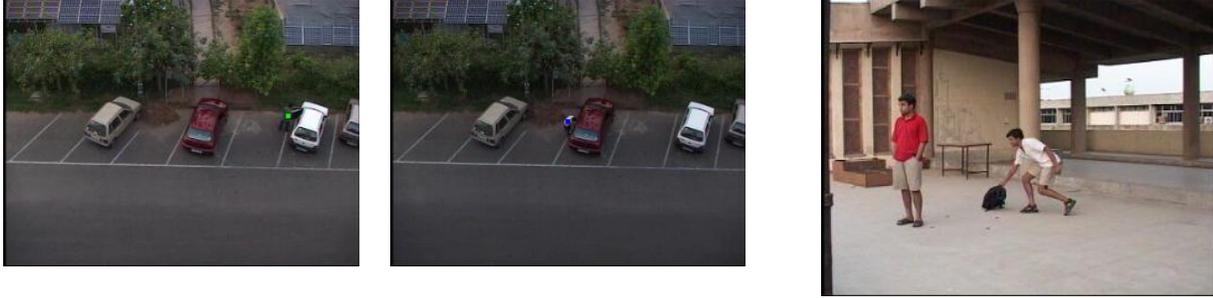


Figure 9: Another example of detection of possible car theft.

- [3] A. R. Chowdhury and R. Chellappa. A Factorization Approach for Activity Recognition. In *CVPR Workshop on Event Mining*, 2003.
- [4] I. Cohen and G. Medioni. Detecting and Tracking Moving Objects for Video Surveillance. In *Proc. Computer Vision and Pattern Recognition*, 1999.
- [5] R. T. Collins and et. al. A system for video surveillance and monitoring. Technical report, CMU-RI-TR-00-12, 2000. Also in Special section on video surveillance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8), August 2000.
- [6] D. Gibbins, G. N. Newsam, and M. Brooks. Detecting Suspicious Background Changes in Video Surveillance of Busy Scenes. In *Proc. of IEEE Workshop on Applications in Computer Vision*, pages 22 – 26, 1996.
- [7] R. Hartley and A. Zisserman. *Multiple view geometry in computer vision*. Cambridge University Press, 2000.
- [8] S. Hongeng, F. Bremond, and R. Nevatia. Representation and Optimal Recognition of Human Activities. In *Proc. Computer Vision and Pattern Recognition*, pages 1818 – 1825, 2000.
- [9] S. Hongeng and R. Nevatia. Multi-agent Event Recognition. In *Proc. Intl. Conf. Computer Vision*, pages 84 – 93, 2001.
- [10] M. V. Hua Zhong, Jianbo Shi. Detecting Unusual Activity in Video. In *Proc. Computer Vision and Pattern Recognition*, page (to appear), 2004.
- [11] Y. Ivanov and A. Bobick. Recognition of Visual Activities and Interactions by Stochastic Parsing. *IEEE Trans. on Pattern Anal. and Machine Intell.*, 22(8):852 – 872, August 2000.
- [12] V. Kettner. Time-dependent HMMs for Visual Intrusion Detection. In *IEEE Workshop on Event Mining: Detection and Recognition of Events in Video*, June 2003.
- [13] A. Madabhushi and J. Aggarwal. A Bayesian Approach to Human Activity Recognition. In *Proc. of the 2nd International Workshop on Visual Surveillance*, pages 25 – 30, 1999.
- [14] G. Medioni, I. Cohen, F. Bremond, S. . Hongeng, and R. Nevatia. Event Detection and Analysis from Video Stream. *IEEE Trans. on Pattern Anal. and Machine Intell.*, 23(8):873 – 889, 2001.
- [15] D. Moore, I. Essa, and M. Hayes. Exploiting Human Actions and Object Context for Recognition Tasks. In *Proc. Intl. Conf. Computer Vision*, pages 80 – 86, 1999.

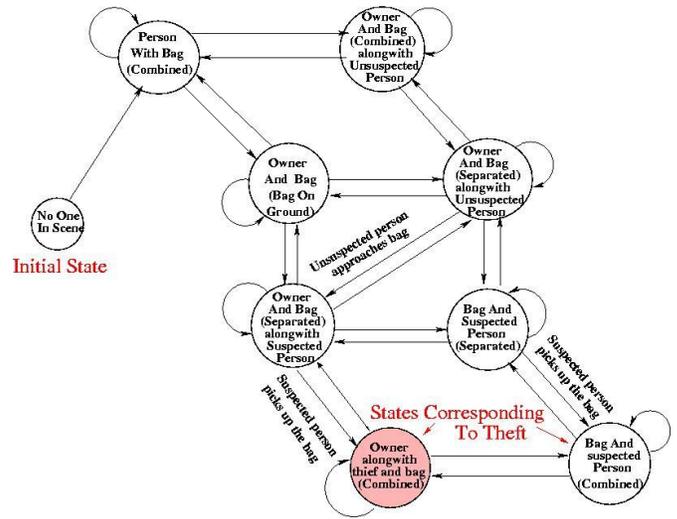


Figure 10: Explicit programming of the stealing example. The event layer FSM is given in the right.

- [16] N. Oliver, E. Horvitz, and A. Garg. Layered Representations for Human Activity Recognition. In *Fourth IEEE Int. Conf. on Multimodal Interfaces*, pages 3 – 8, 2002.
- [17] N. Oliver, B. Rosario, and A. Pentland. A Bayesian Computer Vision System for Modeling Human Interactions. In *Proc. of ICVS*, pages 255 – 272, 1999.
- [18] Y. Shi and A. Bobick. Representation and Recognition of Activity using Propagation Nets. In *Proc. of 16th Intl. Conf. on Vision Interface*, 2003.
- [19] T. Starner and A. Pentland. Visual Recognition of American Sign Language using Hidden Markov Models. In *Intl. Workshop on Automatic Face- and Gesture-Recognition*, 1995.
- [20] C. Stauffer and W. Grimson. Learning Patterns of Activity using Real-time Tracking. *IEEE Trans. on Pattern Anal. and Machine Intell.*, 22(8):747 – 757, 2000.
- [21] N. Vaswani, A. R. Chowdhury, and R. Chellappa. Activity Recognition Using the Dynamics of the Configuration of Interacting Objects. In *Proc. Computer Vision and Pattern Recognition*, pages 633 – 640, 2003.
- [22] J. Yamato, J. Ohya, and K. Ishii. Recognizing Human Action in Time-Sequential Images using Hidden Markov Model. In *Proc. Computer Vision and Pattern Recognition*, pages 379 – 385, 1992.
- [23] L. Zelnik-Manor and M. Irani. Event-based video analysis. In *Proc. Computer Vision and Pattern Recognition*, 2001.