

A VIDEO SYSTEM FOR URBAN SURVEILLANCE: FUNCTION INTEGRATION AND EVALUATION*

R. J. Oliveira

P. C. Ribeiro

J. S. Marques

J. M. Lemos

IST/INESC-INOV

IST/ISR

IST/ISR

IST/INESC-ID

ricardo.oliveira@inov.pt

pribeiro@isr.ist.utl.pt

jsm@isr.ist.utl.pt

jlml@inesc.pt

Abstract: This paper is concerned with video sequence analysis for urban area surveillance applications. The aim is to detect, track and classify targets entering a urban scene under varying illumination conditions and distracters. The paper contributions consist in the integration of algorithms for performing the various tasks and in their statistical evaluation. Results are presented on the basis of a benchmark video sequence.

1. INTRODUCTION

Video sequence analysis for surveillance applications has been the subject of several recent research papers [5, 7, 1]. The system described in most of these works comprise the functions of object detection, tracking, recognition and classification. High level tasks such as object activity recognition has also been addressed, e.g., [5].

The problem of object detection has been tackled using statistical models of the background image [8, 5, 7], frame differences techniques or a combination of both [6]. Several techniques have also been used for object tracking in video sequences in order to cope with multiple interacting targets. These range from Kalman trackers, nearest neighbor trackers [10], to multiple hypothesis trackers such as the PDAF [10], multiple hypothesis tree [4] or long term tracking using Bayesian Networks [1].

Object recognition and classification is performed using statistical Pattern Recognition and neural networks. Several features, which explore the specific condition of the problem can be used. These include geometric features such as bounding box aspect ratio, motion patterns and color histogram [7, 5].

Most of the existing literature address the above problems by considering algorithms for their solution without concern for their evaluation. Evaluation problems have been considered in [9, 3]. These tackle mainly segmentation problems of video processing algorithms aiming at object representation and coding. In [3], ROC curves for performance evaluation of video surveillance processing system are considered. However, only object detection is studied. Furthermore only two types of errors (false alarms and miss detections) are treated, nothing being said about other types of errors, such as region splitting and region

merging, which degrade the performance of the overall system. Therefore, the main contribution of this paper consists in the statistical evaluation of the methods for each of the surveillance tasks. The system considered results from the integration of modules performing the tasks of object detection, classification and recognition.

This paper is organized as follows. After a brief state of the art review which motivates this paper contributions (this section), the surveillance system considered is described in section 2. Experimental evaluation is then performed in section 3. Conclusions are presented in section 4.

2. SYSTEM DESCRIPTION

The surveillance system implemented can be viewed as four independent, but interacting modules: detection, tracking, classification and recognition. The figure 1 describes the system and the interaction between its modules. To per-

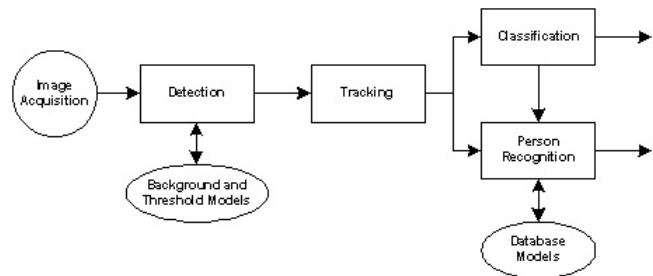


Fig. 1. System Block Diagram.

form the detection task, a robust real-time algorithm, suggested by T. Boulton [8] was adapted. The approach followed uses two adaptive background images, per-pixel adaptive thresholds and a region grouping algorithm, named quasi-connected components(QCC).

The tracking algorithm determines the overlap between detected regions in consecutive frames, in order to link them, when no ambiguity exists. The linking of an active region in consecutive frames originates a stroke, which describes the evolution of the mass center over time.

The classification task is performed each frame for all active regions detected, and the classification of a stroke is performed by determining the most voted class. To cope

* THIS WORK WAS PARTIALLY SUPPORTED BY FEDER AND FCT UNDER PROJECT LTT

with tracking ambiguities, a color-based recognition module is also integrated in the system.

In order to achieve real-time capability, the detection module, which is the most time-consuming one, was implemented in C. This allow the final system to operate over 30 fps with 768x576 images (using an *Intel Centrino 1.3* running *WindowsXP*).

2.1. Detection

Like many systems, our processing starts with a change-detection method based on background subtraction. The fact that undetected targets cannot be tracked, makes detection a crucial stage. The main difficulties of such approach lie in the fact that, even in controlled environments, the background undergoes a continual change, mostly due to the existence of lighting variations and distractors (example: clouds passing by, branches of trees moving with the wind). Target occlusion and interaction with the scene rises additional problems. To overcome these difficulties, the robust and fast algorithm described in [8] was implemented. The robustness towards lighting variations of the scene is achieved using adaptive background models and adaptive per-pixel thresholds. The use of multiple backgrounds and the grouping technique QCC contribute to the robustness of the algorithm towards unwanted distractors.

The system implemented uses two gray scale background models B_1 and B_2 , created during a training phase. The idea is to have both a lower and a higher pixel value, contemplating this way the variations of "non-target" pixels in the scene. The per-pixel threshold, T_L (low threshold), is then initialized to be above the difference between the two backgrounds. A higher threshold T_H is also created, resulting from the addition of a constant value (which represents the sensitivity of the algorithm) to the threshold T_L .

2.2. Tracking

The purpose of tracking is to determine the spatio-temporal information of each target present in the scene. Since the visual motion of targets is always small in comparison to their spatial extends, no position prediction is necessary to construct the strokes [7]. The approach followed associates the active regions of the present frame I^t , with the regions in the previous frame I^{t-1} by region overlap. Five different situations were considered, namely: target entering or leaving the scene, a merge or a split between multiple targets or a match. The association of regions and their classification is based on a binary association matrix computed by testing the overlap of regions in consecutive frames. Whenever there is a match, the stroke is updated. The others situations lead to a stroke interruption (when a target leave), a new stroke (when a target enters) or both (merges and splits).

Tracking also interacts with the detection. When a target stops in the scene for a certain amount of time, the tracker merges the target in the background.

2.3. Classification

For the classification task three main questions must be answered, namely: which classes should be considered, which

features best separate these classes and which classifiers best adapt to the previous choices? One of the main goals of the classifier is to achieve low miss-classification probabilities while considering a wide spectrum of classes. At the same time the goal was not to consider time-dependent features, limiting the classifier exclusively to geometric properties. In this way the resulting classifier can be used in different machines, as it is independent of the achieved frame-rate. The Pets2001 Dataset (Camera1 and Camera2) is set as being a typical working situation, thus the classes considered are: one person, two persons, three persons, one vehicle, two vehicles and mixed groups. The mixed groups consider all other combinations of people and vehicles.

A R^3 feature space was chosen. Two of the chosen characteristics are measurements which take into account the size of the target and the size of its bounding box in pixels.

The third measurement, the *normalized size*, represents how big or small a target is, in comparison with a single person in the same zone \mathcal{Z} of the scene.

$$F_1 = \frac{\text{height}}{\text{width}} \quad F_2 = \frac{\text{target area}}{\text{bounding box area}}$$

$$F_3 = \frac{\text{target area}}{\text{single person area in } \mathcal{Z}}$$

The use of this feature assumes that the relation between the normalized mean size of each class, is constant. Experimental testing revealed that this assumption was acceptable and the overall results were improved about 10%. A training process is necessary to establish the typical size of a person in each zone of the image.

The classes that comprise several merged targets, cannot be described by a gaussian distribution over the feature space. These can assume many different configurations, which makes them harder to parameterize. This suggests the choice of a non-parametric classifier, for example the K-Nearest Neighbors algorithm (KNN), which directly estimates the *a posteriori* probability $P(\omega_i|x)$, using adaptive Parzen Windows.

The classification task interacts with the tracker in each frame, voting for the class of each detected target. In this way, a final class is chosen for each stroke as being the most voted one.

2.4. Recognition

As in the classification module, no time information is used to perform the recognition task. This recognition process is aimed at recognizing in a short term period, i.e. targets that become occluded for a few seconds or targets that merge for a few seconds and then split again. The models are characterized by the *pdf* estimates of the chosen feature space, in this case color. Several histogram constructions under RGB and HSV color space were compared, using both normalized and non-normalized spaces. The best result was achieved using a color histogram in *RGB* color space, comprised of 10 bins of R correlated to 10 bins of G, concatenated with 10 bins of B, giving a total of 110 bins.

The model and the candidate model are represented as follows:

$$\begin{array}{ll} \text{model:} & \hat{\mathbf{q}} = \{\hat{q}_u\}_{u=1\dots m} \quad \sum_{u=1}^m \hat{q}_u = 1 \\ \text{candidate:} & \hat{\mathbf{p}} = \{\hat{p}_u\}_{u=1\dots m} \quad \sum_{u=1}^m \hat{p}_u = 1 \end{array}$$

For the error measure between histograms the Euclidean distance was used.

3. EXPERIMENTAL EVALUATION

To evaluate the performance of the different tasks, the correct results of each task are manually constructed. For the detection task a ground truth was created (Pets2001 Camera1 training and testing sequences) and an automatic method, similar to the binary association matrices, decides if a detection falls into one of the following cases: correct detections, false alarms, miss detections, splits and merges.

For the classification task, the detected targets were manually classified and are distributed by the different classes as follows: 2652 isolated persons, 411 groups of two people, 607 groups of three people, 2998 isolated vehicles, 547 groups of two vehicles and 829 mixed groups.

For isolated persons, seven persons are tracked. A snapshot is chosen as a model for each one, in order to test the recognition methods. The evaluation procedure is done off-line comparing the models of the detected persons with the reference models of all the persons previously seen.

3.1. Working Examples

Below, two working examples of the surveillance algorithm are shown, in which the color represents the class of the objects, as in figure 3. The color of the box around the target represents the voted class in the present frame. The most voted class for each target is represented by the line color and the text below stroke.

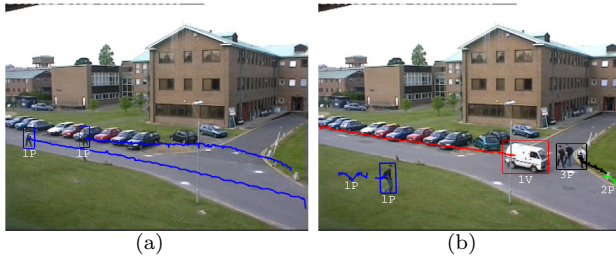


Fig. 2. Working examples

Figure 2b shows several targets being followed simultaneously. The group of three persons was correctly classified as being a group of two at first, because until then, the rightmost element was still out of the scene, and afterwards as a group of three. The path of the leftmost person shows more than one stroke, due to a detection error, in this case a splitting. As explained in section 2.2 these situations make the tracker start new strokes. These strokes can be linked together using the recognition algorithm and some heuristics.

The figure 3 represents the end result of the *Pets2001* sequence analysis, showing all targets path and class.

As groups are the most difficult to classify, figure 4 shows four examples of detected groups and how they enter the feature extraction procedure (2nd row). Due to partial occlusion the rightmost example was misclassified as being a group of two persons.

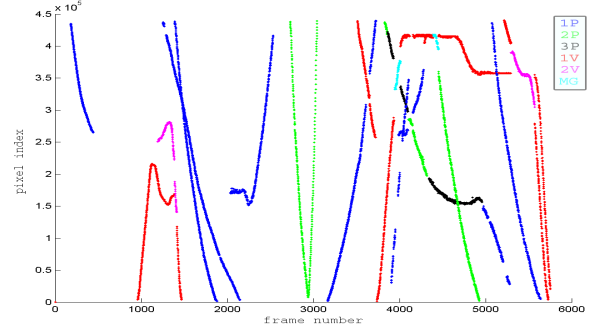


Fig. 3. Stroke map obtained in Pets2001 Camera 1, training and testing sequences



Fig. 4. Detection and classification examples: the rightmost example was incorrectly classified as 2P, the other examples were correctly classified as MG, 3P, MG.

3.2. Statistical Evaluation

The statistical evaluation of a detection algorithm can be done by determining its Receiver Operating Characteristics (ROC) curves [3] (figure 5). The merge or split per frame rate can also be calculated, providing further information about the algorithms characteristics (figure 6). In some applications merge or splits can be very undesirable, for example when used for tracking and classification.

The examples shown in figures 5 and 6 show how the parameter sensitivity defined in chapter 2.1 influences the operation of the detection algorithm. The study of other parameters, further contribute to the use of the full potential of every such algorithm.

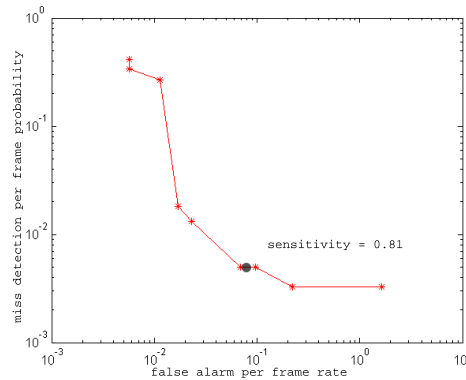


Fig. 5. ROC curve for pairs $(p_{\text{detection}}, r_{\text{fa}})$

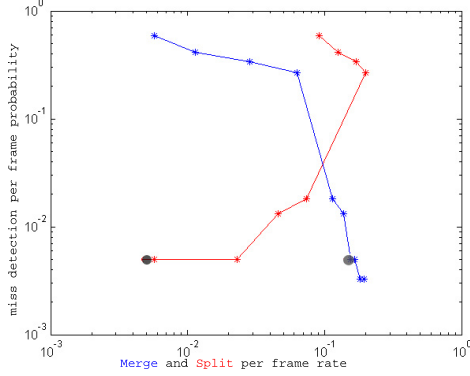


Fig. 6. ROC curve for pairs $(p_{\text{detection}}, r_{\text{merge}})$ and $(p_{\text{detection}}, r_{\text{split}})$.

Target classification and recognition was evaluated using the confusion matrix and the error probability. Assuming a non Gaussian distribution of the data, the K-NN classifier is used. First, the inclusion of a given feature can be tested, to determine if it is relevant to differentiate the chosen classes. By evaluating the inclusion of the feature F_3 (that gives an idea on the size of the target) we reach a 6% and 18% of error probability, with and without F_3 respectively.

The confusion matrix (table 1) show a detailed result of the classifier and help to determine which are the most difficult classes to classify. For example the classes *Mixed Groups* and *Tree Persons* are where it fails the most. These groups include many different configurations of people and vehicles and also many different poses for its elements, leading to a more dispersing feature space.

Table 1. Confusion matrices for the KNN classifier using features: F_1 , F_2 and F_3

	1P	2P	3P	1V	2V	MG
1P	99.3	0.8	0	0	0	0
2P	8.0	87.4	4.1	0.5	0	0
3P	0	3.8	78.3	10.9	0	7.1
1V	0.0	0.2	1.1	96.6	0.8	1.2
2V	0	0	0	4.6	90.1	5.3
MG	0	0	5.9	8.9	4.1	81.1

Although the recognition task is only used in specific occasions: e.g., in the case of entering targets, temporary occlusion or target merge and split, the testing was done in every frame. The table 2 shows a confusion matrix for the seven persons considered in the Pets2001 Camera1 Data set, where an error probability of 18% is achieved. As with the classification algorithm, these tables help determine which features, or in this particular case, which color histogram, best suits this problem.

4. CONCLUSIONS

Evaluation of an integrated real-time video system for urban surveillance has been considered. This system comprises the modules of object detection, classification and recognition. The main contribution of the paper consists in

Table 2. Confusion matrices for RG+B color histograms

	P1	P2	P3	P4	P5	P6	P7
P1	89,0	0	0	0	0	11,0	0
P2	13,1	82,5	4,5	0	0	0	0
P3	0	10,3	89,8	0	0	0	0
P4	5,4	52,3	0	42,3	0	0	0
P5	0	2,7	0	0	97,3	0	0
P6	0	0	0	58,8	0	41,2	0
P7	0	0	0	0	9,4	3,6	87,0

the statistical evaluation of tasks performed by these modules. A benchmark video sequence has been used for the application. The evaluation is based on the comparison of the system output with the ground truth obtained by manually editing the video sequence. Several types of errors were taken into account. The procedure developed can be used as a systematic methodology for video surveillance evaluation.

5. REFERENCES

- [1] A. Abrantes, J. Marques, J. Lemos, "Long Term Tracking Using Bayesian Networks", in *IEEE ICIP*, 609-612, 2002.
- [2] A. Cavallaro, et al. "Objective Evaluation Of Segmented Quality Using Spatio-Temporal Context".in *Proc. of IEEE ICIP*, 22-25 September 2002.
- [3] F. Oberti, et al. "ROC curves for performance evaluation of video sequences processing systems for surveillance applications", DIBE-University of Genoa.
- [4] I. Cox, S. Hingorani, "An Efficient Implementation of Reid's Multiple Hypothesis Tracking Algorithm and its Evaluation for the Purpose of Visual Tracking", *IEEE PAMI*, 138-150, 1996.
- [5] I. Hariataoglu, et al. "W4: A Real-Time Surveillance of People and Their Activities", *IEEE PAMI*, 22, No. 8, August 2000.
- [6] R. Collins, et al. "A System for Video Surveillance and Monitoring", CMU-RI-TR-00-12, 2000.
- [7] S. McKenna, et al, "Tracking Groups of People", *CVIU* 80, 42-56 (2000).
- [8] T. Boulton, et al. "Into the Woods: Visual Surveillance of Noncooperative and Camouflaged Targets in Complex outdoor Settings", in *Proceedings of the IEEE*, vol 89, NO 10 October 2001.
- [9] Y. Zhang, "A Survey on Evaluation Methods for Image Segmentation", *Pattern Recognition*, Vol 29. No. 8, pp.1335-1346, 1996.
- [10] Y Bar-Shalom, T. Fortmann, "Tracking and Data Association", Academic Press, 1988.

Acknowledgement: The authors would like to thank Dr. Jacinto Nascimento for helpful suggestions concerning error types.