

Vehicle and Pedestrian Detection in eSafety Applications

S. Álvarez, M. A. Sotelo, I. Parra, D. F. Llorca, M. Gavilán *

Abstract—This paper describes a target detection system on road environments based on Support Vector Machine (SVM) and monocular vision. The final goal is to provide pedestrian-to-car and car-to-car time gap. The challenge is to use a single camera as input, in order to achieve a low cost final system that meets the requirements needed to undertake serial production in automotive industry. The basic feature of the detected objects are first located in the image using vision and then combined with a SVM-based classifier. An intelligent learning approach is proposed in order to better deal with objects variability, illumination conditions, partial occlusions and rotations. A large database containing thousands of object examples has been created for learning purposes. The classifier is trained using SVM in order to be able to classify pedestrians, cars and trucks. In the paper, we present and discuss the results achieved up to date in real traffic conditions.

Keywords: *Vision, pedestrian detection, vehicle detection, SVM (Support Vector Machine) and tracking.*

1 Introduction

This paper describes a vision-based candidate extraction method for pedestrian and vehicle detection in Intelligent Transportation Systems (ITS). These methods are a challenging problem in real traffic scenarios since they must perform robustly under variable illumination conditions, variable rotated positions and pose, and even if some of the object parts are partially occluded. An additional difficulty is given by the fact that the camera is installed on a fast-moving vehicle. As a consequence of this, the background is no longer static, and candidates significantly vary in scale. This makes the ITS problem quite different from that of detecting and tracking objects in the context of surveillance applications, where the cameras are fixed and the background is stationary.

To ease the recognition task in vision-based systems, a candidate selection mechanism is usually applied. The selection of candidates can be implemented by performing an object segmentation in either a 3-D scene or a 2-D image plane.

*Department of Electronics, University of Alcalá, Madrid (Spain). Email: {sergio.alvarez,sotelo,parra,llorca,miguel.gavilan}@depeca.uah.es

Not many authors have tackled the problem of monocular pedestrian [1],[2] and vehicle recognition [10]. Concerning the various approaches proposed in the literature, most of them are based on shape analysis. On the one hand, about pedestrian detection, some authors use feature-based techniques, such as recognition by vertical linear features, symmetry, and human templates [2], Haar wavelet representation [5], hierarchical shape templates on Chamfer distance [3], correlation with probabilistic human templates [6], sparse Gabor filters and support vector machines (SVMs) [7], and principal component analysis [8].

On the other hand, some previous developments use available sensing methods for vehicle detection, such as radar [14], stereo vision or a combination of stereo-vision and laser [12]. In [9] the authors propose the fusion of stereo-vision and radar for creating a hybrid velocity adaptive control system called HACC. Only a few works deal with the problem of monocular vehicle detection using symmetry and colour features [10] or pattern recognition techniques [13]. In [10] the authors propose the use of horizontal edges and vertical symmetry together with a shape-dependent process for removing objects that are too small or too big in the image plane. In [11] the authors propose the use of a geometrical model for vehicle characterization using evolutionary algorithms, assigning different geometrical models depending on the vehicle lane.

The remaining of the paper is organized as follows: Section 2 provides a description of the pedestrian candidate selection mechanism. Section 3 describes the vehicle detection method. The implementation and results achieved up to date are presented and discussed in Section 4. And finally, Section 5 summarizes the conclusions.

2 Pedestrian Detection

To detect the different pedestrians on the image, two methods are used, depending on their kind of movement and their distance from the camera: *motion analysis* and *texture analysis*. Motion is the main extraction method to select candidates who are moving, because it detects each movement in the image; but there are complex situations which need further exhaustive analysis. The idea is to select the candidates by following the next rule: in-

ward movement with motion analysis, standing candidates with texture analysis and outward movement by a decision between both methods depending on the situation. After selecting the candidates, a classification and tracking step is done, in order to decrease the number of false positive detections and keep the valid candidates under detection.

2.1 Inward Motion

This algorithm is an adaptation of the crowd detection system of [1], which uses the temporal consistency of the gray levels in the image to detect the regions with high probability of being moving. The gray level information of selected regions of the image in consecutive frames is used to estimate the motion on that region. The motion system architecture loops through the following modules: *xt image*, *turning detection*, *inward motion detection* and *candidate generation*.

Xt image

Firstly, the gray intensity level from 19 different horizontal lines in the image (exploration lines) is stored for the 16 last frames. This will be called the *xt image* (Figure 1(b)). Next, the gradient is computed for the *xt image*. Movement will be represented in the *xt gradient image* by lines with different inclinations, depending on the orientation and velocity of the movement (1(c)).

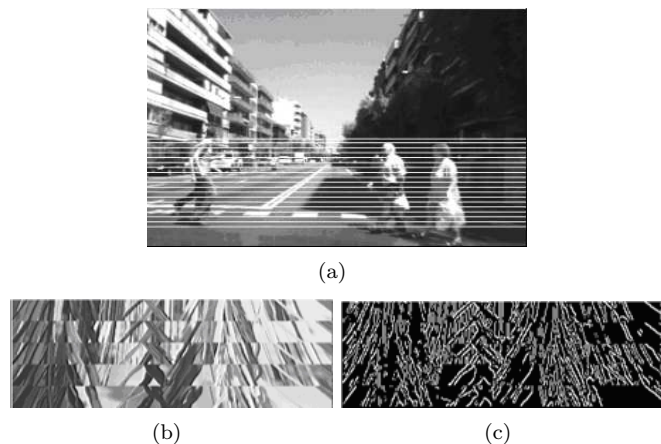


Figure 1: Scan-lines (a) and 5 example lines of processing (b and c) from the *xt image*.

A moving object will form a line in the *xt image*, which will be detected using Hough transform.

Turning detection

This module tracks a sparse set of points in the image and determines if they all move in one direction (figure 2).

The inward motion detection module updates its thresholds according to the output of the “turning detection” module.

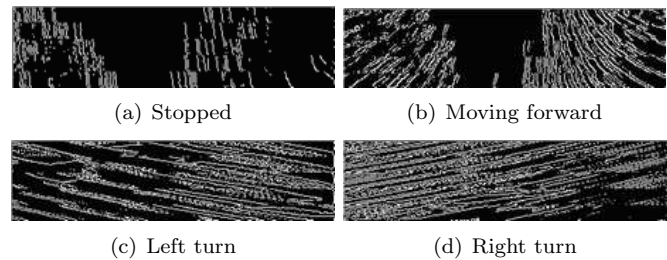


Figure 2: *xt gradient images* for the turning detection module.

Inward motion detection

Analyzing the *xt image* with the Hough-like transform, a probability distribution function is computed for each exploration line. Then inward motion is detected if the probability of left motion in the right part of the image, or vice versa, is above a predefined adaptive threshold. In Figure 3, movement for the main scan-lines in (a) is shown (b).

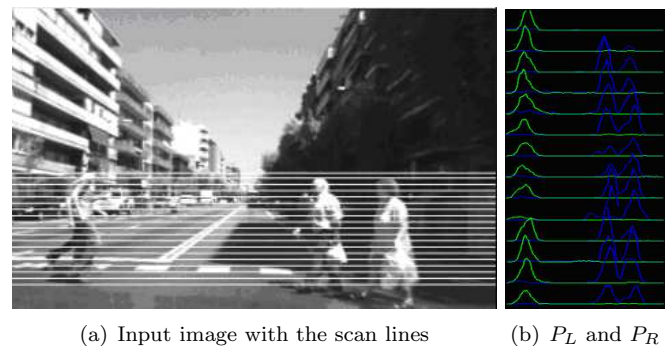


Figure 3: Left and right motion for the scan lines.

Candidates generation

Using the projective integral for the scan-lines along the *x*-axis a global movement image can be obtained. The probability of a column of having left movement is given by (and similarly, for the right motion):

$$p_L(x) = \sum_N L(x, n) \tag{1}$$

where *x*, varying from 0 to the image width, represents the column and *N* is the number of exploration lines.

This image will give an idea of the “amount of movement per column” at every pixel of the exploration lines (see

figure 4). This approximation allows for an easy and fast candidate extraction based on “columns of movement”. Given the camera position and the average pedestrian size, it is impossible to have two separate pedestrians in the image plane that contribute to the same column of movement. In the case of overlapping pedestrians contributing to the same column of movement the closest one will be tracked which is the system expected behaviour.

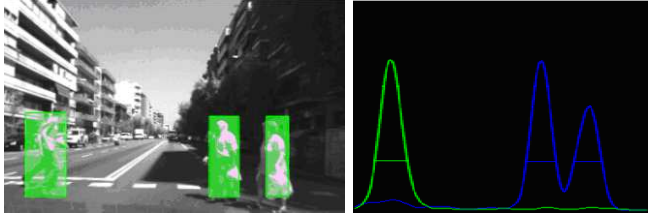


Figure 4: Example of 2D candidates and its movement.

Adaptive thresholds are set to determine the candidate columns depending on the input from the turning detection system and on the local average motion. The local regions on which the adaptive thresholds are applied are usually fixed to three, although this can be modified by the input from the turning detection system to adapt the system to hard turns.

Every column above its threshold will create a “candidate column” indicating that movement was detected for that column. This will set the candidates width. Next, every scan-line will be checked against the created candidates to decide if this line belongs to the new candidate. If movement was detected on this scan-line for the candidate column then this scan line will be part of the candidate. This will set the candidates height. The resulting motion candidates will be bounding boxes on the image.

Finally, a last consideration is taken. Because of the vehicle movement, pedestrian velocity increases, and as the algorithm works with the previous 16 frames, the system can lose spatial resolution. To avoid this situation, an image scaling is done depending on the information of the car velocity.

2.2 Texture Analysis

The main target of this system will be the pedestrians moving laterally and towards the path of the vehicle. The ego-motion of the vehicle and the cluttered environments make the motion maps very noisy at the sides when moving forward or at the front when turning. In these situations it is not easy to separate the movements of the candidates from the movement of the background due to the ego-motion. In fact, when the vehicle is moving the velocities of a longitudinally walking pedestrian and of a wall, street lamp, door, etc, behind them are very similar. A different detection technique is thus needed for these scenarios: texture analysis.

The texture analysis consists of two basic steps: *Image preprocessing*, where image textures are extracted, protecting vertical edges, an important pedestrian characteristic; and *Spatial localization*, to find the different candidates of the image based on the high entropy areas.

Image preprocessing

To reduce image noise, an appropriate method is to filter the captured image using median filters to remove small errors (figure 5(a)). Relevant information is then extracted with a texture analysis. In this case, a calculation of the local range of the image is done (rangefilt), which provides information about the local variability of the pixel intensity levels as figure 5(b) shows. As can be seen, uniform areas are set to zero and the lines of texture change are set to a value higher than zero. Then, a border extraction is done by providing output pixels which contains the range value (maximum value - minimum value) of the 3x3 neighbourhood around the corresponding pixel in the input image. The main advantage of this method (above other edges extractors like canny or sobel for example), is the lack of thresholds, so illumination changes are less important in this algorithm.



(a) Smoothed image

(b) Rangefilt image

Figure 5: Texture analysis example.

After this extraction, a thresholding is done in order to binarize the image, and finally, an erode step with a vertical kernel is used to preserve the vertical edges.

Spatial Localization

After preprocessing, a Region Of Interest (ROI) with the ideal dimensions of a pedestrian (width and height) is sequentially projected along the image in order to find areas with information that could belong to a pedestrian. These ideal dimensions are obtained, for each vertical coordinate, by a perspective camera model called pinhole [16], and the coordinate system shown in figure 6.

After obtaining the projected dimensions of a pedestrian for each vertical coordinate, the image is sequentially scanned from the bottom to the horizon line, looking for ROIs with enough entropy to be a valid candidate (more than an experimental threshold). Entropy is computed

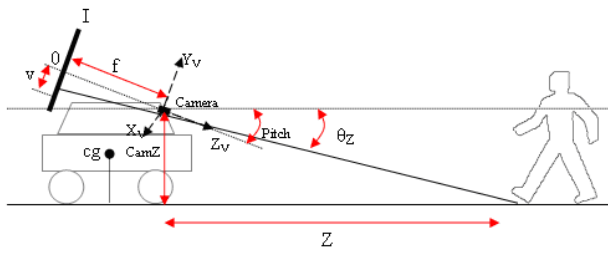


Figure 6: Vertical road mapping geometry.

as the ratio of connected vertical edge points normalized by the size of the area being analysed.

If the projected ROIs do not have enough entropy, or if it is distributed vertically only in a part of the area, the candidate is discarded (figure 7(a)). Otherwise this coefficient is analyzed to detect if it has central distribution. Then the candidates selected will have enough entropy (distributed vertically) and centred horizontal distribution of this, as figure 7(c) shows.

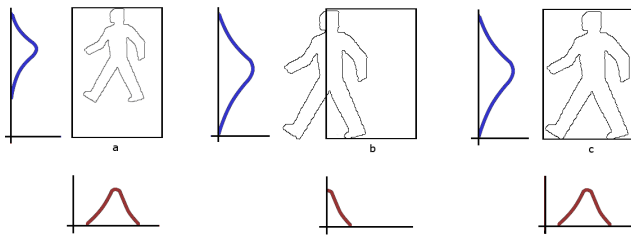


Figure 7: Different distributions of the entropy.

In a last step, candidate region overlapping is computed and regions inside other regions or with large overlapping are discarded. Figure 8 shows a real example about the textures method used. As can be seen, standing pedestrians are detected (two pedestrians on the left), and the problem with turning situations is solved (pedestrian on the right); so the combination of the two explained methods cover most of the possible situations.



Figure 8: Candidates extraction using texture analysis.

2.3 Pedestrian Recognition and Tracking

Pedestrian recognition is carried out by using a learning-based approach in which discriminative features are extracted from each candidate and then, they are passed

through the learning machine or classifier (SVM classifier).

The tracking process relies on Kalman filter theory to provide spatial estimates of detected pedestrians. Thus, steadier 2D spatial positions of pedestrian in the image plane are obtained. In addition, Kalman filter provides the prediction of the 2D position of each candidate in the next frame, which will be very useful when solving the data association problem, that is, the association of the candidates state vectors at frame i and the measurements available at frame $i + 1$.

The overall structure of this part of the algorithm is depicted in Figure 9.

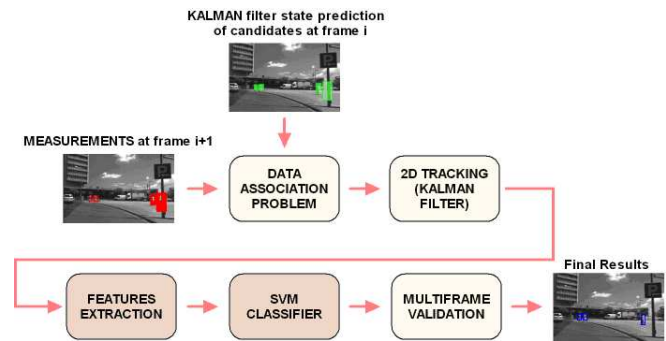


Figure 9: Overall scheme of tracking and recognition steps.

3 Vehicle Detection

3.1 Candidate Selection

The main idea of the vehicle candidate selection is to make a sequential scanner for each road lane, from the bottom to the horizon line of the image looking for collections of horizontal edges that might represent a potential vehicle. The scanned lines are associated in groups of three. For each group, a horizontality coefficient is computed as the ratio of connected horizontal edge points normalized by the size of the area being analysed. The resulting coefficient is used together with a symmetry analysis in order to trigger the attention mechanism.

An adaptive thresholding process is implemented in order to obtain robust edges from the road images. This adaptive process is based on an iterative algorithm that gradually increases the contrast of the image, and compares the number of edges obtained in the contrast increased image with the number of edges obtained in the actual image. If the number of edges in the actual image is higher than in the contrast increased image the algorithm stops. Otherwise, the contrast is gradually increased and the process resumed. After thresholding, horizontal edges in the scanned regions given by a Lane Departure Warning System (LDWS), developed by the authors in previous works [15], are examined to detect the rear part of potential vehicles. In order to decide if the collection of

horizontal lines represents a possible vehicle candidate, its width is compared to that of an ideal car. The ideal car width is obtained for each vertical coordinate using the camera pinhole model explained in section 2.2.

Once the car width is computed at the current frame it is compared to the collection of horizontal lines found after the thresholding analysis. If they are similar to some extent defined by an empirical value, a square area above the collection of horizontal lines, denoted as candidate ROI, is considered for further analysis. The aim is to compute the entropy of the candidate ROI and its vertical symmetry. Only those regions containing enough entropy and symmetry are identified as potential vehicles. Figure 10 shows a detailed block diagram of the detection procedure and figure 11(a) depicts an example of the detection step.

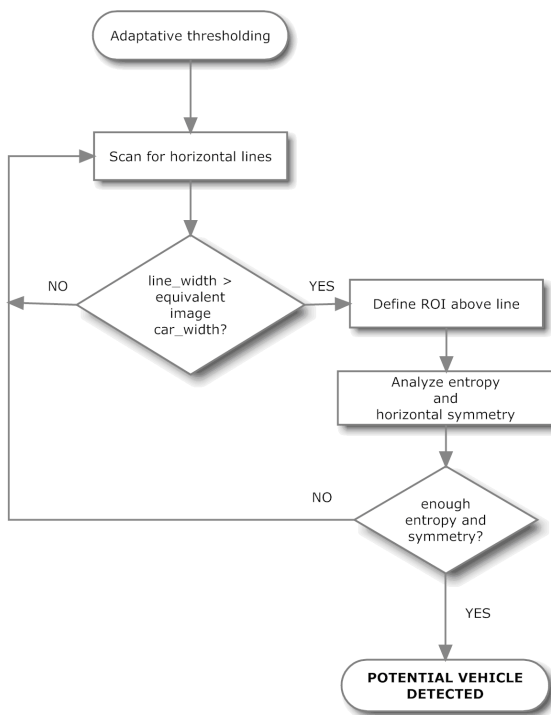


Figure 10: Block diagram of the vehicle detection mechanism.

Similarly to pedestrian detection, detected candidates are classified as vehicles or non-vehicles depending on features obtained from the vehicle ROI using Support Vector Machines (SVM), and after that they are tracked using Kalman filter techniques. Figure 11(b) shows the result of the classification step.

The purpose of the Kalman filtering is to obtain a more stable position of the detected vehicles. Besides, oscillations in vehicles position due to the unevenness of the road makes y coordinate of the detected vehicles change several pixels up or down. This effect makes the distance detection unstable, so a Kalman filter is necessary for

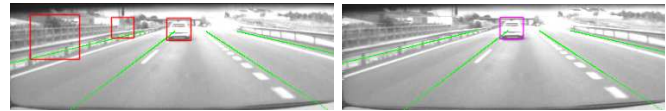


Figure 11: Result of the detection step (left) and after SVM classification (right)

minimizing these kinds of oscillations.

4 Implementation and Results

The system was implemented on a Pentium IV PC at 2.4 GHz, onboard a Citroën C4, running the Knoppix GNU/Linux Operating System and Libsvm libraries. Using 320 x 240 pixel images, the complete algorithm runs at an average rate of 20 frames/s, depending on the number of pedestrians and vehicles being tracked and their position. The candidate selection system has proved to be robust in various illumination conditions, different scenes, and distances up to 25m for pedestrians and 40m for vehicles.

The size of the databases created to generate the SVM models (67.650 samples for pedestrian and 10.000 samples for vehicles) represents a crucial factor to take care of. To obtain a sufficiently representative set we have taken pedestrians and cars and trucks as positive samples; and lamppost, litter bins, crash barriers, median strip, pieces of road, etc, like negative samples; all of them taken in different weather conditions (with and without rain, shadows, etc).

The quality of the classification system is mainly measured by means of the detection rate/false positive ratio (DR/FPR). These two indicators are graphically bounded together in an ROC (figure 12).

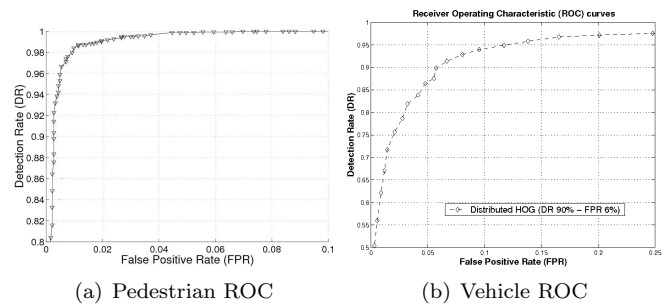


Figure 12: Receiver Operating Characteristics (ROC Curve).

The selection of the FPR value has been made to show performance in representative points where differences between curves can be optimally appreciated. FPR must be a value for which DR exhibits an acceptable value. This leads to selecting DR of 99% and FPR of 2% for pedestrians and DR of 90% and FPR of 5% for vehicles.

5 Conclusions

We have developed and implemented a pedestrian & vehicle detection system based on Support Vector Machine (SVM) and monocular vision with the objective of providing pedestrian-to-car and car-to-car time gap measurement for applications in the framework of Intelligent Transportation Systems (ITS). Candidates are raised using an attention mechanism based on inward motion, textures, horizontal edges, vertical symmetry and entropy. The detected objects are passed on to a SVM-based classifier. After classification, detected objects are tracked using Kalman filtering. A large database containing thousands of examples extracted from real images has been created for learning purposes. After assessment of the practical results achieved in our experiments, the following general conclusions can be summarized:

- Car dynamics (yaw rate and velocity) have to be taken into account in order to improve data association and tracking stages. Moreover, at present, the region of interest is statically fixed. By using yaw rate and car velocity variables we can define a more precise region of interest and evaluate the risk for each candidate.
- The performance of the vehicle detection module is significantly increased by building on the output provided by the LDWS function.
- The presence of large shadows on the asphalt due to vehicles circulating along the road produces negative effects on the candidate selection mechanism, yielding to inaccuracy in measuring the distance to the vehicles.

Acknowledgements

This work has been supported by means of Research Grant P9/08 from the Spanish Ministry of Development.

References

- [1] A. Shashua, Y. Gdalyahu and G. Hayun, "Pedestrian detection for driving assistance systems: Single-frame classification and system level performance", *Proc. IEEE Intell. Veh. Symp.*, Parma, Italy, pp. 1-6, 2004.
- [2] A. Broggi, M. Bertozzi, A. Fascioli and M. Sechi, "Shape-based pedestrian detection", *Proc. IEEE Intell. Veh. Symp.*, Dearborn, MI, Oct. 2000.
- [3] D. M. Gavrila and V. Philomin, "Real-time object detection for smart vehicles", *Proc. 7th IEEE Int. Conf. Comput. Vis.*, 1999, pp. 8793.
- [4] I. Parra, D. Fernández, M. A. Sotelo, "Combination of Feature Extraction Methods for SVM Pedestrian Detection", *IEEE Trans. Int. Transp. Sys.*, 2007. vol 8. no 2.
- [5] C. Papageorgiou and T. Poggio, "A trainable system for object detection," *Int. J. Comput. Vis.*, vol. 38, no. 1, pp. 1533, 2000.
- [6] H. Nanda and L. Davis, "Probabilistic template based pedestrian detection in infrared videos," in *Proc. IEEE Intell. Veh. Symp.*, Versailles, France, Jun. 2002, pp. 1520.
- [7] H. Cheng, N. Zheng, and J. Qin, "Pedestrian detection using sparse Gabor filter and support vector machine," *Proc. IEEE Intell. Veh. Symp.*, Las Vegas, NV, Jun. 2005, pp. 583587.
- [8] U. Franke, D. Gavrila, S. Gorzig, F. Lindner, F. Puetzold, and C. Wohler, "Autonomous driving goes downtown," *IEEE Intell. Syst. Their Appl.*, vol. 13, no. 6, pp. 4048, Nov./Dec. 1998.
- [9] L. Bombini, P. Cerri, P. Medici, and G. Alessandretti. "Radar-vision fusion for vehicle detection". In *Procs. International Workshop on Intelligent Transportation*, pp. 65-70, Hamburg, Germany, 2006.
- [10] A. Broggi, P. Cerri, and P. C. Antonello. "Multi-Resolution Vehicle Detection using Artificial Vision". *IEEE Intelligent Vehicles Symposium*. Parma, Italy. 2004.
- [11] C. Hilario, J. M. Collado, J. M. Armingol, and A. de la Escalera. "Visual Perception and Tracking of Vehicles for Driver Assistance Systems". *IEEE Intelligent Vehicles Symposium*. 2006.
- [12] R. Labayrade, C. Royere, D. Gruyer, and D. Aubert, "Cooperative fusion for multi-obstacles detection with use of stereovision and laser scanner". In *Proc. Int. Conf. On Advanced Robotics*, 2003.
- [13] G. P. Stein, O. Mano, and A. Shashua, "Vision-based ACC with a single camera: bounds on range and range rate accuracy". In *Proc. Int. Conf. Intelligent Vehicles*, 2002.
- [14] G. R. Widman, W. A. Bauson, and S. W. Alland, "Development of collision avoidance systems at Delphi Automotive Systems". In *Proc. Int. Conf. Intelligent Vehicles*, pp. 353-358, 1998.
- [15] M. A. Sotelo, J. Nuevo, L. M. Bergasa, M. Ocaña, "A monocular solution to vision-based ACC in road vehicles", *Lecture Notes in Computer Science.*, Vol. 3643, pp. 507-512 2005.
- [16] P.F. Alcantarilla, L.M. Bergasa, P. Jiménez, "Night Time Vehicle Detection for Driving Assistance Light-Beam Controller". *IEEE Intelligent Vehicles Symposium*. 2008.