

Video analysis based vehicle detection and tracking using an MCMC sampling framework

Jon Arróspide¹, Luis Salgado¹, Marcos Nieto²

¹ Image Processing Group (GTI), Polytechnic University of Madrid, Madrid, 28040, Spain (phone: +34-91-3367353; e-mail: {jal, lsa}@gti.ssr.upm.es)

² Vicomtech-IK4, Research Alliance, San Sebastián, 20009, Spain (phone: +34-943-309230; e-mail: mnieto@vicomtech.org)

Received: 15 May 2011 / Revised version: date

Abstract This paper presents a probabilistic method for vehicle detection and tracking through the analysis of monocular images obtained from a vehicle-mounted camera. The method is designed to address the main shortcomings of traditional particle filtering approaches for use in traffic environments. Namely, Bayesian methods based on importance sampling do not scale well as the dimensionality of the feature space grows, which involves important limitations when it comes to tracking of multiple objects. Alternatively, the proposed method is based on a Markov chain Monte Carlo (MCMC) approach, which allows efficient sampling of the feature space. The method involves important contributions as regards both the motion

and the observation models of the tracker. Indeed, as opposed to particle filter-based tracking methods in the literature, which typically resort to observation models based on appearance or template matching, in this work a likelihood model that combines appearance analysis with information from motion parallax is introduced. Regarding the motion model, a new interaction treatment is defined based on Markov Random Fields (MRF) that allows to handle possible inter-dependencies in vehicle trajectories. As for vehicle detection, the method relies on a supervised classification stage using Support Vector Machines (SVM). The contribution in this field is two-fold. First, the mostly rectilinear structure of vehicles is capitalized on to define a new descriptor based on the analysis of gradient orientations in concentric rectangles. This descriptor involves a much smaller feature space compared to traditional descriptors, which are too costly for real-time applications. Second, a new vehicle image database is generated to train the SVM and made public. The proposed vehicle detection and tracking method is proven to outperform existing methods and to successfully cope with the challenging situations contained in the test sequences.

Key words Image processing, object tracking, Monte Carlo methods, intelligent vehicles

1 Introduction

Signal processing techniques have been widely used in sensing applications to automatically characterize the environment and for understanding of the scene. Typical problems include ego-motion estimation, obstacle detection, or object localization, monitoring and tracking, which are usually addressed by processing of the information coming from sensors such as Radar, LIDAR, GPS or video-cameras. Specifically, methods based on video analysis play an outstanding role due to their low cost, the striking increase on processing capabilities, and the significant advances in the field of computer vision.

Naturally object localization and monitoring are crucial to a good understanding of the scene. However, they are especially critical in safety applications where the objects may constitute a threat to the observer or to any other individual. In particular, tracking of vehicles in traffic scenarios from an on-board camera constitutes a major focus of scientific and commercial interest, as vehicles originate the majority of accidents.

Video-based vehicle detection and tracking have been addressed in many different ways in the literature. The former aims at localizing vehicles by exhaustive search in the images, whereas the latter pursues to keep track of already detected vehicles. As regards vehicle detection, since exhaustive search throughout the image is costly, most of the methods in the literature proceed in a two-stage fashion: hypothesis generation, and hypothesis verification. The former usually involves a rapid search so that the image regions

that do not match some expected feature of the vehicle are disregarded and only a small number regions potentially containing vehicles are further analyzed. Typical features include edges [1], color [2,3], and shadows [4]. Many works based on stereovision have also been proposed (e.g. [5,6]), although they involve a number of drawbacks with respect to monocular methods, especially in terms of cost and flexibility.

Verification of hypotheses is usually addressed through model-based or appearance-based techniques. Model-based techniques exploit the a priori known structure of the vehicles to generate a description (i.e., the model) that can be matched with the hypotheses to decide whether it is a vehicle or not. Both rigid (e.g. [7]) and deformable (e.g. [8]) vehicle models have been proposed. Appearance-based techniques, in contrast, involve a training stage in which features are extracted from a set of positive and negative samples to design a classifier. Neural Networks [9] and Support Vector Machines [10, 11] are extensively used for classification, while many different possibilities have been proposed for feature extraction. Among others, Histograms of Oriented Gradients (HOG) [12,13], Principal Component Analysis [14], Gabor filters [11] and Haar-like features [15,16] have been applied to derive the feature set for classification.

Direct use of many of these techniques is very time-consuming and thus unrealistic in real-time applications. Therefore, in this paper we propose a vehicle detection method that exploits the intrinsic structure of the vehicles in order to achieve good detection results while involving a small feature

space (and hence low computational overhead). The method combines prior knowledge on the structure of the vehicle, based on the analysis of vertical symmetry of its rear, with appearance-based feature training using a new HOG-based descriptor and SVM. Additionally, a new database containing vehicle and non-vehicle images has been generated and made public, which is used to train the classifier. The database separates between vehicle instances depending on their relative position with respect to the camera and hence allows to adapt the feature selection and the classifier in the training phase according to the vehicle pose.

As regards object tracking, feature-based and model-based approaches have been traditionally utilized. The former aim at characterizing the objects by a set of features (e.g., corners [17] and edges [18] have been used to represent vehicle) and to subsequently track the object through inter-frame feature matching. In contrast, model-based tracking uses a template that represents a typical instance of the object, which is often dynamically updated [19,20]. Unfortunately, both approaches are prone to errors in traffic environments due to the difficulty to extract reliable features or to provide a canonical pattern of the vehicle.

To confront these problems, many recent approaches to object tracking entail a probabilistic framework. In particular, the Bayesian approach [21, 22], especially in the form of particle filtering, has been used in many recent works (e.g. [23–25]), to model the inherent degree of uncertainty in the information obtained from image analysis. Bayesian tracking of multiple ob-

jects can be found in the literature both using individual Kalman or Particle Filters (PF) for each object [26,24], and a joint filter for all the objects [27, 28]. The latter is better suited for applications in which there is some degree of interaction between objects, as it allows to control the relations between objects in a common dynamic model (those are much more complicated to handle through individual particle filters [29]). Notwithstanding, the computational complexity of joint-state traditional importance sampling strategies grows exponentially with the number of objects, which results in a degraded performance with respect to independent PF-based tracking when there are several participants (as occurs in the traffic scenario).

On the other hand, PF-based object tracking methods found in the literature resort to appearance information for the definition of the observation model. For instance, in [23], a likelihood model comprising edge and silhouette observation is employed to track the motion of humans. In turn, the appearance-based model used in [27] for ant tracking consists of simple intensity templates. However, methods using appearance-only models are only bound to be successful under controlled scenarios, such as those in which the background is static. In contrast, the considered on-board traffic monitoring scenarios entail a dynamically changing background and varying illumination conditions, which affect the appearance of the vehicles. Hence, observation models based only on appearance are prone to errors.

In this work we present a new framework for vehicle tracking which combines efficient sampling, handling of vehicle interaction, and reliable obser-

vation modeling. The proposed method is based on the use of Markov chain Monte Carlo approach to sampling (instead of the traditional importance sampling) which renders joint state modeling of the objects affordable, while also allowing to easily accommodate interaction modeling. In effect, driver decisions are affected by neighboring vehicle trajectories (vehicle tends to occupy free space), therefore an interaction model based on Markov Random Fields [30] is introduced to manage inter-vehicle relations. In addition, an enriched observation model is proposed, which fuses appearance information with motion information. Indeed, motion is an inherent feature of vehicles and is exploited here through the geometric analysis of the scene. Specifically, the projective transformation relating the road plane between consecutive time points is instantaneously derived and filtered temporally based on a data estimation framework using a Kalman filter. The difference between the current image and the previous image warped with this projectivity allows to detect regions likely featuring motion. Most importantly, the combination of appearance and motion based information provides robust tracking even if one of the sources is temporarily unreliable or not available. The proposed system has proven to successfully track vehicles in a wide variety of challenging driving situations and to outperform existing methods.

This paper is organized as follows. In Section 2 the reader is introduced to the problem of Bayesian tracking, and the general framework of the method proposed to address it is presented. While in Section 2 the con-

tributions on each of the constituent parts of the framework are hinted, in-depth descriptions of those are provided in the following sections. In particular, Section 3 describes the details of the vehicle tracking algorithm based on MCMC sampling. In turn, Section 4 presents the designed motion model including interaction treatment, whereas Section 5 and 6 address the definition of the the observation model regarding appearance-based analysis and motion-based analysis, respectively. The description of the method proposed for vehicle detection as well as details regarding the database used to train it are enclosed in Section 7. Section 8 comprises the experiments conducted on the vehicle detection and tracking method and a discussion on the observed results, followed by the conclusions at the end of the paper.

2 Overview of the proposed framework

As explained in the introduction, the proposed tracking method is grounded on a Bayesian inference framework. Object tracking is addressed as a recursive state estimation problem in which the state consists of the positions of the objects. The Bayesian approach allows to recursively update the state of the system upon receipt of new measurements. If we denote \mathbf{s}_k the state of the system at time k and \mathbf{z}_k the measurement at the same instant, then Bayesian theory provides an optimal solution for the posterior distribution of the state given by

$$p(\mathbf{s}_k|\mathbf{z}_{1:k}) = \frac{p(\mathbf{z}_k|\mathbf{s}_k) \int p(\mathbf{s}_k|\mathbf{s}_{k-1})p(\mathbf{s}_{k-1}|\mathbf{z}_{1:k-1})d\mathbf{s}_{k-1}}{p(\mathbf{z}_k|\mathbf{z}_{1:k-1})} \quad (1)$$

where $\mathbf{z}_{1:k}$ integrates all the measurements up to time k [21]. Unfortunately, the analytical solution is intractable except for a set of restrictive cases. Particularly, when the state sequence evolution is a known linear process with Gaussian noise and the measurement is a known linear function of the state (also with Gaussian noise) then the Kalman filter constitutes the optimal algorithm to solve the Bayesian tracking problem. However, these conditions are highly restrictive and do not hold for many practical applications. Hence, a number of suboptimal algorithms have been developed to approximate the analytical solution. Among them, particles filters (also known as bootstrap filtering or condensation algorithm) play an outstanding role and have been used extensively to solve problems of very different nature. The key idea of particles filters is to represent the posterior probability density function by a set of random discrete samples (called particles). In the most common approach to particles filtering, known as importance sampling, the samples are drawn independently from a proposal distribution $q(\cdot)$, called importance density. In addition, each sample is assigned a weight which depends on its likelihood $p(\mathbf{z}_k|\mathbf{s}_k)$.

However, importance sampling is not the only approach to particle filtering. In particular, Markov chain Monte Carlo methods provide an alternative framework in which the particles are generated sequentially in a Markov chain using at each step the approximation to the posterior distribution. In this case, all the samples are equally weighted and the solution in (1) can therefore be approximated as

$$p(\mathbf{s}_k | \mathbf{z}_{1:k}) \approx c \cdot p(\mathbf{z}_k | \mathbf{s}_k) \sum_{r=1}^N p(\mathbf{s}_k | \mathbf{s}_{k-1}^{(r)}) \quad (2)$$

where the state of the r -th particle at time k is denoted $\mathbf{s}_k^{(r)}$, N is the number of particles, and c is the inverse of the evidence factor in the denominator of (1). The advantage of MCMC methods is that the complexity increases only linearly with the number of objects, in contrast to importance sampling, in which the complexity grows exponentially [27]. This implies that using the same computational resources, MCMC will be able to generate a larger number of particles and hence to better approximate the posterior distribution than importance sampling. Therefore, in this work an MCMC framework is adopted for vehicle tracking. The general scheme of the method is summarized in Fig. 1.

This framework requires definition of the observation model, $p(\mathbf{z}_k | \mathbf{s}_k)$, and the dynamic or motion model, $p(\mathbf{s}_k | \mathbf{s}_{k-1})$. The motion model is designed under the assumption that vehicles velocity can be approximated to be locally constant, which is valid in highway environments. As a result, the evolution of a vehicle's position can be traced by a first-order linear model. However, linearity is lost due to the perspective effect in the acquired image sequence. To preserve linearity we resort to a plane rectification technique, usually known as Inverse Perspective Mapping (IPM) [31]. This computes the projective transformation, T , that produces an aerial or bird's-eye view of the scene from the original image. The image resulting of plane rectification will be referred to as the rectified domain or the transformed domain.

In the rectified domain, the motion of vehicles can be safely described as a first-order linear equation with an added random noise.

One important issue regarding the dynamic model is the interaction treatment. Most approaches to multiple vehicle tracking involve an independent motion model for each vehicle. However, this requires some external method for handling of interaction, and often this is simply disregarded. In contrast, we have designed an MRF-based interaction model that can be easily integrated with the above-mentioned individual vehicle dynamic model.

On the other hand, the observation model is critical to the performance of the method, hence much effort is devoted to its design. Likelihood models are typically built according to the observation of the appearance of the objects. Here, we extend the likelihood model so that it not only includes a set of appearance-based features but also considers a feature that is inherent to the vehicles, i.e., their motion. In particular, the model for the observation of motion is based on the temporal alignment of the images in the sequence through the analysis of multiple-view geometry. As for the appearance-based observation model, rather than usual template matching methods, a probabilistic approach is defined using a Expectation Maximization approach for likelihood function optimization.

Finally, a method is necessary to detect new vehicles in the scene so that these can be integrated in the tracking framework. This is addressed in the current work by using a two-step procedure composed of an initial

hypothesis generation and a subsequent hypothesis verification. In particular, candidates are verified using a supervised classification strategy over a new descriptor based on HOG features. The proposed feature descriptor and the classification strategy are explained in Section 7.

3 Vehicle tracking algorithm

The designed vehicle tracking algorithm aims at estimating the position of the vehicles existing at each time of the image sequence. Hence, the state vector is defined to comprise the position of all the vehicles $\mathbf{s}_k = \{s_{i,k}\}_{i=1}^M$, where $s_{i,k}$ denotes the position of vehicle i , and M is the number of vehicles existing in the image at time k . As stated, the position of a vehicle is defined in the rectified domain given by the transformation T , although back-projection to the original domain is naturally possible via the inverse projective transformation T^{-1} .

An example of the bird's-eye view obtained through IPM is illustrated in Fig. 2. Observe that the upper part of the vehicles is distorted in the rectified domain. This is due to the fact that IPM calculates the appropriate transformation for a given reference plane (in this case the road plane), which is not valid for all the elements outside this plane. Therefore, analysis is focused on the road plane and the position of a vehicle will be defined as the middle point of its lower edge, i.e., the contact point between the road and the vehicle.

In order to estimate the joint state of all the vehicles, the MCMC method is exploited. As mentioned, in MCMC the approximation to the posterior distribution of the state is given by (2), which assuming that the likelihood of the different objects is independent can be rewritten as follows:

$$p(\mathbf{s}_k | \mathbf{z}_{1:k}) \approx c \cdot \prod_{i=1}^M p(z_{i,k} | s_{i,k}) \sum_{r=1}^N p(\mathbf{s}_k | \mathbf{s}_{k-1}^{(r)}) \quad (3)$$

where $z_{i,k}$ is the observation at time k for object i . The Markov chain of samples at time k is generated as follows. First, the initial state is obtained as the mean of the samples in $k - 1$, $\mathbf{s}_k^0 = \sum_r \mathbf{s}_{k-1}^{(r)} / N$. New samples for the chain are generated from a proposal distribution $Q(\cdot)$. In particular, we follow a Gibbs-like approach, in which only one target is changed at each step of the chain. At step τ the proposed position $s'_{i,k}$ of the randomly picked target i is thus sampled from the proposal distribution, which in our case is a Gaussian centered at the value of the last sample for that target, $Q(s'_{i,k} | s_{i,k}^{(\tau-1)}) = \mathcal{N}(s'_{i,k} | s_{i,k}^{(\tau-1)}, \sigma_q)$. The candidate sample is therefore $\mathbf{s}'_{i,k} = (\mathbf{s}_{\setminus i,k}^{(\tau-1)}, s'_{i,k})$, where $\mathbf{s}_{\setminus i,k}$ denotes \mathbf{s}_k but with $s_{i,k}$ omitted. This sample is accepted or not according to the Metropolis algorithm, which evaluates the posterior probability of the candidate sample in comparison to that of the previous sample and defines the following probability of acceptance [30]:

$$A(\mathbf{s}'_k, \mathbf{s}_k^{(\tau-1)}) = \min \left(1, \frac{p(\mathbf{s}'_k | \mathbf{z}_{1:k})}{p(\mathbf{s}_k^{(\tau-1)} | \mathbf{z}_{1:k})} \right) \quad (4)$$

This implies that if the posterior probability of the candidate sample is larger than that of $\mathbf{s}_k^{(\tau-1)}$ the candidate sample is accepted, and if it is

smaller, it is accepted with probability equal to the ratio between them. In the case of acceptance, $\mathbf{s}_k^{(\tau)} = \mathbf{s}'_k$. Otherwise the previous sample is replicated $\mathbf{s}_k^{(\tau)} = \mathbf{s}_k^{(\tau-1)}$.

Observe that the samples obtained with the explained procedure are highly correlated. It is a common practice to retain only every L -th sample and leave out the rest, which is called thin-out. In addition, the first B samples are discarded to prevent the estimation from being degraded by bad initialization. Finally, at each time step the vehicle position estimates, $\bar{\mathbf{s}}_k = \{\bar{s}_{i,k}\}_{i=1}^M$, are inferred as the mean of the valid particles $\mathbf{s}_k^{(r)}$:

$$\bar{\mathbf{s}}_k = \frac{1}{N} \sum_{r=1}^N \mathbf{s}_k^{(r)} \quad (5)$$

4 Motion and interaction model

The motion model is defined in two steps: the first layer copes with the individual movement of a vehicle in the absence of other participants, and the second layer addresses the movement of vehicles in a common space. As stated in Section 2, the motion of vehicles in the rectified domain is modeled to be linear with constant velocity. Hence, the dynamic model for an individual vehicle can be synthesized by the following equation:

$$s_{i,k} = s_{i,k-1} + v\Delta t + m_k \quad (6)$$

where v is the velocity of the vehicle, which is estimated from previous time points, Δt is the elapsed time between frames, and m_k is an i.i.d. Gaussian noise sequence. The individual dynamic model can be reformulated as

$$p(s_{i,k}|s_{i,k-1}) = \mathcal{N}(s_{i,k}|s_{i,k-1} + v\Delta t, \sigma_m) \quad (7)$$

where σ_m is the variance vector of m_k .

Once the expected evolution of each individual target has been defined, their interaction must also be accounted for in the model. A commonly used way to address interaction is through MRFs (Markov Random Fields), which graphically represent a set of conditional independent relations. An MRF (also known as undirected graph) is composed of a set of nodes V , which represent the variables, and a set of links representing the relations between them. The joint distribution of the variables can be factorized as a product of functions defined over subsets of connected nodes (called cliques, \mathbf{x}_C). These functions are known as potential functions and denoted $\phi_C(\mathbf{x}_C)$. In the proposed MRF the nodes V_i (representing the vehicle positions $s_{i,k} = \{x_{i,k}, y_{i,k}\}$) are connected according to a distance-based criterion. Specifically, if two vehicles, i and j , are at a distance smaller than a predefined threshold, then the nodes representing the vehicles are connected and form a clique. The potential function of the clique is defined as

$$\phi_C(\mathbf{x}_C) = 1 - \exp\left(-\frac{\alpha_x \delta x^2}{w_l^2}\right) \exp\left(-\frac{\alpha_y \delta y^2}{d_s^2}\right) \quad (8)$$

where $\delta x = |x_{i,k} - x_{j,k}|$ and $\delta y = |y_{i,k} - y_{j,k}|$. The functions $\phi_C(\mathbf{x}_C)$ can be regarded as penalization factors that decrease the joint probability of a hypothesized state if it involves unexpected relations between targets. Potential functions consider the expected width of the lane, w_l , and the longitudinal safety distance, d_s . In addition, the design parameters α_x and α_y are selected so that $\alpha_x = 0.5$ and $\alpha_y = 0.5$ whenever a vehicle is at a distance $\delta x = w_l/4$ or $\delta y = d_s$ of another vehicle. Finally, the joint probability is given by the product of the individual probabilities associated to each node and the product of potential functions in existing cliques:

$$p(\mathbf{s}_k | \mathbf{s}_{k-1}) = \prod_{i=1}^M p(s_{i,k} | s_{i,k-1}) \prod_{\mathcal{C}} \phi_C(\mathbf{x}_C) \quad (9)$$

where \mathcal{C} is the set of the two-node cliques. Let us now introduce this motion model in the expression of the posterior distribution in (2):

$$p(\mathbf{s}_k | \mathbf{z}_{1:k}) \approx c \cdot p(\mathbf{z}_k | \mathbf{s}_k) \sum_{r=1}^N \prod_{i=1}^M p(s_{i,k} | s_{i,k-1}^{(r)}) \prod_{\mathcal{C}} \phi_C(\mathbf{x}_C) \quad (10)$$

It is important to note that the potential factor does not depend on the previous state, therefore (10) can be rewritten as

$$p(\mathbf{s}_k | \mathbf{z}_{1:k}) \approx c \cdot p(\mathbf{z}_k | \mathbf{s}_k) \prod_{\mathcal{C}} \phi_C(\mathbf{x}_C) \sum_{r=1}^N \prod_{i=1}^M p(s_{i,k} | s_{i,k-1}^{(r)}) \quad (11)$$

Modeling of vehicle interaction thus requires only the evaluation of an additional factor in the posterior approximation, while producing significant gain in the tracking performance.

5 Appearance-based analysis

The first part of the observation model deals with the appearance of the objects. The aim is to obtain the probability $p_a(z_{i,k}|s_{i,k})$ of the current appearance observation given the object state $s_{i,k}$ (note the subscript a that denotes "appearance"). In other words we would like to know if the current appearance-related measurements support the hypothesized object state. In order to derive the probability $p_a(z_{i,k}|s_{i,k})$ we will proceed in two levels. First, the probability that a pixel belongs to a vehicle will be defined according to the observation for that pixel. Second, by analyzing the pixel-wise information around the position given by $s_{i,k}$, the final observation model will be defined at region level.

The pixel-wise model aims at providing the probability that a pixel belongs to a vehicle. This will be addressed as a classification problem, and it is therefore necessary to define the different categories expected in the image. In particular, the rectified image (see example in Fig. 2) contains mainly three types of elements: vehicles, road pavement and lane markings. A fourth class will also be included in the model to account for any other kind of elements (such as median stripes or guard rails).

The Bayesian approach is adopted to address this classification problem. Specifically, the four classes are denoted by $\mathcal{S} = \{P, L, V, U\}$, which correspond to the pavement, lane markings, vehicles and unidentified elements. Let us also denote X_i the event that a pixel x is classified as belonging to the class $i \in \mathcal{S}$. Then, if the current measurement for pixel x is represented

by z_x , the posterior probability that the pixel x corresponds to X_i is given by the Bayes rule

$$P(X_i|z_x) = \frac{p(z_x|X_i)P(X_i)}{P(z_x)} \quad (12)$$

where $p(z_x|X_i)$ is the likelihood function, $P(X_i)$ is the prior probability of class X_i , and $P(z_x)$ is the evidence, computed as $P(z_x) = \sum_{i \in \mathcal{S}} p(z_x|X_i)P(X_i)$, which is a scale factor that ensures that the posterior probabilities sum to one. The likelihoods and prior probabilities are defined in the following section.

5.1 Likelihood functions

In order to construct the likelihood functions, a set of features have to be defined that constitute the current observation regarding appearance. These features should achieve a good degree of separation between classes, while at the same time being significant for a broad set of scenarios. In general terms the following considerations hold when analyzing the appearance of the bird's-eye view images. First, the road pavement is usually homogeneous with slight intensity variations among pixels. In turn, lane markings constitute near-vertical stripes of high-intensity, surrounded by regions of lower intensity. As for vehicles, they typically feature very low intensity regions in their lower part, due to vehicle's shadow and wheels. Hence, two features are used for the definition of the appearance-based likelihood model, namely the intensity value, I_x , and the response to a lane-marking detector, I_x . For

the latter, any of the methods available in the literature can be utilized [31, 32]. For this work, a lane marking detector similar to that presented in [33] is used, whose response is defined in every row of the image as

$$R_x = 2I_x - (I_{x-\tau} + I_{x+\tau}) \quad (13)$$

where τ is the expected width of a lane marking in the rectified domain. The likelihood models are defined as parametric functions of these two features. In particular, they are modeled as Gaussian probability density functions:

$$p(I_x|X_i) = \frac{1}{\sqrt{2\pi}\sigma_{I,i}} \exp\left(-\frac{1}{2\sigma_{I,i}^2}(I_x - \mu_{I,i})^2\right) \quad (14)$$

$$p(R_x|X_i) = \frac{1}{\sqrt{2\pi}\sigma_{R,i}} \exp\left(-\frac{1}{2\sigma_{R,i}^2}(R_x - \mu_{R,i})^2\right) \quad (15)$$

where the parameters for the intensity and the lane marking detector are denoted respectively by the subscripts ‘I’ and ‘R’. Specifically, the distribution corresponding to the unknown class, which would intuitively be uniformly distributed for both features, is instead also modeled to be a Gaussian of very high fixed variance to ease further processing (the gain will be clear in Section 5.1.1). All these parameters will be estimated by means of an optimization process based on the Expectation-Maximization (EM) algorithm, explained below. Additionally, likelihood functions are assumed to be conditionally independent on these features for all the classes X_i , thus it is

$$p(z_x|X_i) = p(I_x|X_i)p(R_x|X_i) \quad (16)$$

5.1.1 Parameter estimation: As mentioned, the parameters of the likelihood models in (14) and (15) are estimated via EM. This technique enables us to obtain the maximum likelihood estimate of the parameters of a distribution from a set of observed data. The data distribution is given in this case by

$$p(I_x) = \sum_{i \in \mathcal{S}} p(X_i)p(I_x|X_i) \quad (17)$$

$$p(R_x) = \sum_{i \in \mathcal{S}} p(X_i)p(R_x|X_i) \quad (18)$$

Since the densities of the features I_x and R_x are independent, the optimization is carried out separately for these features. Let us first rewrite the expression (17) so that the dependence on the parameters is explicit:

$$p(I_x|\Theta_I) = \sum_{i \in \mathcal{S}} \omega_{I,i} p(I_x|\Theta_{I,i}) \quad (19)$$

where $\Theta_{I,i} = \{\mu_{I,i}, \sigma_{I,i}\}$ and $\Theta_I = \{\Theta_{I,i}\}_{i \in P,L,V}$. Observe that the prior probabilities have been substituted by factors $\omega_{I,i}$ to adopt the notation typical of mixture models. In effect, EM provides an analytical solution to Gaussian mixture-density parameter estimation problems, such as the one posed here. This has been extensively exploited in many problems, as explained in [34], where the analytical solutions can also be found. In this

kind of problems, the set of unknown parameters is composed of the parameters of the densities and of the mixing coefficients, $\Theta = \{\Theta_{I,i}, \omega_{I,i}\}_{i \in P,L,V}$. Thereby, the parameters resulting from the final EM iteration are fed into the Bayesian model defined in equations (12)-(15). The process is completely analogous for the feature R_x .

The EM algorithm is proven to converge to a local maximum, hence it is necessary to provide a good starting point. In this case, since EM is applied in every frame of the incoming image sequence, the results from the previous image can recursively be used as starting point in the current frame. However, initialization is still necessary for the triggering of the process. This is addressed through the analysis of histograms of each of the features in the bird's-eye view image. As regards the intensity feature, first the pixels that are likely to correspond to the pavement class are selected by filtering out the regions of the image that feature high gradient. This is done through the application of the Sobel operator, followed by thresholding and morphological dilation. An example of resulting binary mask is illustrated in Fig. 3 (b) for the image in (a). The image obtained after high-gradient pixel removal is shown in Fig. 3 (c), in which we see that only the pavement pixels are retained. A histogram is then generated from the remaining pixels (see Fig. 4 (a)), and the initial parameters for $\mu_{I,P}$ and $\sigma_{I,P}$ are extracted from it.

The corresponding histograms for the lane marking and the vehicle class are generated in turn by taking the values above $\mu_{I,P} + \sigma_{I,P}$ and below

$\mu_{I,P} - \sigma_{I,P}$, respectively. This satisfies the preliminary assumption that $\mu_{I,O} < \mu_{I,P} < \mu_{I,L}$. The corresponding map of pixels are shown in Fig. 3 (d) and (e) for the example image in Fig. 3 (a). We proceed in an analogous way to that explained for the pavement class to extract the initial values of $\{\mu_{I,L}, \sigma_{I,L}\}$ and $\{\mu_{I,V}, \sigma_{I,V}\}$. The corresponding histograms are shown in Fig. 4 (b)-(c). A Gaussian distribution with the inferred parameters is also printed on the histograms.

As for the second feature, i.e. the response to the lane marking filter, the histogram typically comprises two modes. The first mode, located in the low R values, corresponds to the homogeneous road and to the regions occupied by vehicles. The second mode covers a broad range of values corresponding to the lane markings (which are imaged with varying intensity depending on the illumination and the distance). By finding the appropriate threshold, histograms can be derived for the two modes, and the corresponding parameters are used for initialization. Note that this feature does not discriminate between vehicles and pavement, thus it is $\{\mu_{R,P}, \sigma_{R,P}\} = \{\mu_{R,V}, \sigma_{R,V}\}$.

5.2 Appearance-based likelihood model

The result of the proposed appearance-based likelihood model is a set of pixel-wise probabilities of each of the classes. Naturally, in order to know the likelihood of the current object state candidate, we must evaluate the region around the vehicle position given by $s_{i,k} = (x_{i,k}, y_{i,k})$. The vehicle position has been defined as the midpoint of its lower edge (i.e., the seg-

ment delimiting the transition from road to vehicle). Hence, we expect that in the neighborhood above $s_{i,k}$, pixels display high probability to belong to the vehicle class, $p(X_V|x)$, while the neighborhood below $s_{i,k}$ should involve low vehicle probabilities if the candidate state is good. Therefore, the appearance-based likelihood of the object state $s_{i,k}$ is defined as

$$p_a(z_{i,k}|s_{i,k}) = \frac{1}{(w+1)h} \left(\sum_{x \in R_a} p(X_V|z_x) + \sum_{x \in R_b} (1 - p(X_V|z_x)) \right) \quad (20)$$

where R_a is the region of size $(w+1) \times h/2$ above $s_{i,k}$, $R_a = \{x_{i,k} - w/2 \leq x < x_{i,k} + w/2; y_{i,k} - h/2 \leq y < y_{i,k}\}$, and R_b is the region below $s_{i,k}$, $R_b = \{x_{i,k} - w/2 \leq x < x_{i,k} + w/2; y_{i,k} < y \leq y_{i,k} + h/2\}$.

6 Motion-based analysis

As mentioned above, the second source of information for the definition of the likelihood model is motion analysis. Two-view geometry fundamentals are used to relate the previous and current views of the scene. In particular, the homography (i.e. projective transformation) of the road plane is estimated between these two points in time. This allows us to generate a prediction of the road plane appearance in future instants. However, vehicles (which are generally the only objects moving on the road plane) feature inherent motion in time, hence their projected position in the plane differs from that observed. The regions involving motion are identified through image alignment of the current image and the previous image warped with the homography. These regions will correspond to vehicles with high probability.

6.1 Homography calculation

The first step towards image alignment is the calculation of the road plane homography between consecutive frames. As shown in [35] this can be obtained from a minimum of four feature correspondences by means of the Direct Linear Transformation (DLT). In this work features are extracted through the Harris detector [36] and matched using KLT [37], although any other standard technique, such as SIFT [38], can be used for this purpose.

Although a homography estimate H is now available from the DLT applied over the correspondences, straightforward image alignment is not possible. Indeed, it must be taken into account that this homography might be highly unreliable due to the following reasons. First, the road is usually homogeneous and thus the number of features resulting from standard feature extraction techniques is small. In addition, the inclusion of inaccurate or wrong correspondences for the computation of DLT is especially harmful when the number of points is small.

Therefore, intermediate processing of the computed homography is necessary. This is achieved in the present work by means of a linear estimation process based on Kalman filtering. Let us first inspect the analytical expression of the homography between two consecutive instants. Fig. 5 illustrates the situation of a vehicle with an on-board camera moving on a flat road plane, $\pi_0 = (\mathbf{n}^\top, d)^\top$, where $\mathbf{n} = (0, 1, 0)^\top$ and d is the distance between the camera and the ground plane. The coordinate system of the camera at time k_1 is adopted as the world coordinate system. At time k_2 the camera

has moved to position \mathbf{C}_2 , and rotation $R_x(\alpha)$ might have occurred around X-axis due to camera shaking. Additional rotation $R_y(\beta)$ in the Y-axis must be considered in the case the vehicle changes lane or takes a curve. From the previous discussion, and assuming a pinhole camera model, the camera projection matrices at times k_1 and k_2 are respectively

$$\begin{aligned} P_1 &= K[I|\mathbf{0}] \\ P_2 &= KR_x(\alpha)R_y(\beta)[I - \mathbf{C}_2] \end{aligned} \quad (21)$$

The homography H relates the projections, \mathbf{x}_1 and \mathbf{x}_2 , of a 3D point, $\mathbf{X} \in \pi_0$, in two different images. Its expression can be derived from the equations in (21). In effect, for the first view it is $\mathbf{x}_1 = P_1\mathbf{X} = K[I|\mathbf{0}]$ and hence any point in the ray $\mathbf{X} = (\mathbf{x}_1^\top(K^{-1})^\top, \rho)^\top$ projects to \mathbf{x}_1 . The intersection of this ray and the plane π_0 determines the value of the parameter ρ : it is $\pi_0^\top\mathbf{X} = \mathbf{n}^\top K^{-1}\mathbf{x} + d\rho = 0$, and thus $\rho = -\mathbf{n}^\top K^{-1}\mathbf{x}_1/d$. The projection of the point \mathbf{X} into the second view is given by

$$\begin{aligned} \mathbf{x}_2 &= P_2\mathbf{X} = KR_x(\alpha)R_y(\beta)[I - \mathbf{C}_2]\mathbf{X} = \\ &= KR_x(\alpha)[R_y(\beta)| - R_y(\beta)\mathbf{C}_2](\mathbf{x}_1^\top(K^{-1})^\top, \rho)^\top = \\ &= KR_x(\alpha)[R_y(\beta)K^{-1}\mathbf{x}_1 + \mathbf{t}\rho] = \\ &= KR_x(\alpha)[R_y(\beta) - \mathbf{t}\mathbf{n}^\top/d]K^{-1}\mathbf{x}_1 \end{aligned}$$

where $\mathbf{t} = -R_y(\beta)\mathbf{C}_2$. This vector constitutes the translation in the direction of heading of the vehicle and is thus given by $\mathbf{t} = (0, 0, 1)^\top v/f_r$, where

v is the velocity of the vehicle and f_r is the frame rate. From the above equations the expression of the homography of the plane π_0 between k_1 and k_2 is derived:

$$\mathbf{H} = \mathbf{K}\mathbf{R}_x(\alpha)[\mathbf{R}_y(\beta) - \mathbf{t}\mathbf{n}^\top/d]\mathbf{K}^{-1} \quad (22)$$

6.1.1 Time-filtering framework: At each time k we have a noisy approximation of the homography \mathbf{H} of the road plane between the previous and the current instant. However, the evolution of \mathbf{H} in time is assumed to be smooth due to the intrinsic constraints in the vehicle dynamics, therefore better estimates can be obtained by filtering noisy measurements in time. Temporally filtered estimates of the homography are obtained by modeling \mathbf{H} with a zero-order Kalman filter whose state vector is composed of the elements H_{ij} of the homography matrix. The design of the filter is summarized as follows:

$$\mathbf{x}_k^\top = \{H_{ij}, 1 \leq i, j \leq 3\}$$

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \mathbf{w}_k$$

$$\mathbf{z}_k^\top = \{H_{ij}^k, 1 \leq i, j \leq 3\}$$

$$\mathbf{z}_k = \mathbf{x}_k + \mathbf{v}_k$$

The process and measurement noise, \mathbf{w}_k and \mathbf{v}_k , are assumed to be given by independent Gaussian distributions, $p(w) \sim N(0, \mathbf{Q})$ and $p(v) \sim N(0, \mathbf{R})$.

Observe that the measurement vector is composed of the elements of the

instantaneous homography matrix, H^k , computed from image correspondences. As stated above, measurements are expected to be prone to error due to the usually small set of correspondences available, hence the measurement error should be tuned to be larger than the process noise (in the proposed configuration it is $Q = 10^{-6}$, $R = 10^{-3}$).

The designed filter provides corrected estimates for the homography at time k , \hat{H}^k , built from the posterior estimate of the filter state, $\hat{\mathbf{x}}_k$. Most importantly, this measure can be used as a prediction for the homography in the next time point. This prediction provides an effective reference to evaluate whether the computed instantaneous measurement may be erroneous or not. Indeed, at the current time k , we can compare the instantaneous homography H^k to the prediction made in the previous time instant \hat{H}^{k-1} : if H^k is close to the expected value \hat{H}^{k-1} then the filter equations will be conveniently updated; in contrast, if the matrices are significantly different, then it is natural to think that noisy correspondences were involved in the calculation of H^k .

The distance between matrices is measured according to the norm of the matrix of differences. Specifically, the norm induced by the 2-norm of a Euclidean space is used. This is obtained by performing Singular Value Decomposition (SVD) of the matrix and retaining its largest singular value [39]. The incoming matrices are accepted and introduced into the Kalman filtering framework only if $\|H^k - \hat{H}^{k-1}\| < t_a$. Otherwise, the measured homography is deemed to be unreliable and the predicted homography is used.

The threshold t_a modulates the maximum acceptable distance to the predicted matrix, which depends on the kinematic restrictions of the platform in which the camera is mounted.

In the case of highways, vehicle dynamics are bounded by the maximum speed, the maximum turning angle (i.e., yaw angle, β) and the maximum variation in the pitch angle, α , for a given frame rate. The maximum velocity is considered to be $v = 120$ km/h (33.3 m/s), as enforced by most nation governments. Besides, a maximum pitch angle variation of $\alpha = \pm 5^\circ$ is considered, and an upper bound of $\beta = \pm 3^\circ$ is inferred for the turning angle according to the standard road geometry design rules. Taking into account these bounds, and assuming an image processing rate of at least 1 fps, the threshold is experimentally found to be $t_a = 60$.

6.2 Motion-based likelihood model

Once a time-filtered estimate of the homography \hat{H}^k is available, reliable image alignment can be performed. Image alignment allows to locate the regions of the image likely featuring motion (and therefore likely containing vehicles). The previous image is aligned with the current image by warping it with \hat{H}^k . Image alignment is exemplified in Fig. 6. In the upper row the snapshots of a sequence at times $k - 1$ are k and displayed. In Fig. 6 (c), the image in (a) warped with \hat{H}^k is shown. Observe that this is very similar in the road region to the actual image at time k (Fig. 6 (b)).

As suggested in the overview of Section 6, the reason for image alignment is that all elements in the road plane (except for the points of the vehicle that belong to this plane) are static, and thus their actual position matches that projected by the homography. In contrast, vehicles are moving, hence their positions in the road plane at time k significantly differs from that projected by the homography, which assumes they are static. Therefore, the differences between the image at time k and the image at time $k - 1$ warped with \hat{H}^k shall be null for all the elements of the road plane except for the contact zones of the vehicles with the road. The differences in these regions will be more significant the larger the velocity of the vehicles. Fig. 6 (d) illustrates the difference between the current image -Fig. 6 (b)- and the previous image warped -Fig. 6 (c)- for the example referred below. As can be observed, whiter pixels -indicating significant difference- appear in the areas of motion of the vehicles in the road. The transformation of the elements outside the road is naturally not well-represented by \hat{H}^k (this is the homography of the road plane) and thus random regions of high differences arise in the background, which will be considered as clutter.

The pixel-wise difference between the current image and the previous image warped provides information on the likelihood of the current object state candidate, $s_{i,k}$. Analogously to the appearance-based likelihood modeling, the region around the vehicle position indicated by $s_{i,k}$ will be evaluated in order to derive its likelihood. Also, to preserve the duality with the appearance-based analysis, the processing is shifted to the recti-

fied domain using the transformation T defined in Section 2. The resulting image, denoted D_r , is illustrated in Fig. 6 (e) for the previous example. In particular, the likelihood of belonging to a region of motion is maximum in $x_{max} = \text{argmax}(D_r(x))$, hence a map of probabilities that the pixel x belongs to a moving region, denoted $p(m|x)$, can be inferred for the whole image as $p(m|x) = D_r(x)/D_r(x_{max})$.

As follows from the above discussion, observe that the regions of high difference are between the current vehicle position and the position that it would occupy if it were static (which is always closer to the camera). Therefore, as opposed to the appearance-based modeling (Section 5.2), we expect that in the neighborhood below the current vehicle position, $s_{i,k}$, pixels have high likelihood values $p(m|x)$, whereas the neighborhood above x should involve small or null probabilities of motion. Hence, the likelihood of the current vehicle state $s_{i,k}$ regarding the motion analysis is defined as

$$p_m(z_{i,k}|s_{i,k}) = \frac{1}{(w+1)h} \left(\sum_{x \in R_a} (1 - p(m|x)) + \sum_{x \in R_b} p(m|x) \right) \quad (23)$$

where the regions R_a and R_b are those defined in Section 5.2, and the subscript m in the probability denotes that it refers to motion observation. The likelihood result obtained from the motion-based analysis is finally combined with that achieved after appearance-based analysis. The joint likelihood of a candidate state $s_{i,k}$ is simply defined as the arithmetic mean of likelihoods:

$$p(z_{i,k}|s_{i,k}) = \frac{1}{2} (p_a(z_{i,k}|s_{i,k}) + p_m(z_{i,k}|s_{i,k})) \quad (24)$$

Note that, although the product of likelihoods could have been used instead, the mean is preferred in order to avoid that the calculation is biased by too small likelihood values.

7 Vehicle detection

Up to this point the method for vehicle tracking has been explained. However, in normal driving situations it is natural that vehicles come in and out of the field of view of the camera throughout the sequence of images. While management of outgoing vehicles is fairly straightforward (the track simply exceeds the limits of the image), a method for incoming vehicles must be designed. The method proposed in this work follows a two-step approach. In the first stage, hypotheses for vehicle positions are made using the results of appearance-based classification explained in Section 5. In the second, those are verified according to the analysis of a set of features in their associated regions in the original domain.

7.1 Hypothesis generation

Exhaustive search of a certain pattern in the whole image is too time-consuming for applications requiring real-time operation. Hence, it is usual to perform some kind of fast pre-processing that restricts the search areas. In this case, we exploit the information extracted from the construction

of likelihood models for tracking and use it to generate a set of candidate regions that will be further analyzed. In particular, two types of inputs could be used corresponding to the appearance-based analysis in Section 5 and the motion-based analysis in Section 6. As referred in the corresponding section, the latter usually involves noise due to background structures, thus appearance-based information is more suitable for hypothesis generation.

Specifically, based on the appearance analysis, for each pixel the probability that it belongs to a vehicle, $p(X_V|z_x)$ is available. We expect that if there is a new vehicle appearing in the image a compact zone of high probabilities must be observed in the surrounding of its position. Therefore, in order to localize new vehicles, a binary map B_m is created containing the pixels in which the probability of the vehicle class is larger than that of the other classes, $p(X_V|\mathbf{z}_x) > p(X_i|\mathbf{z}_x)$, $i \in P, L, U$. As an example, the binary map obtained for the image in Fig. 7 (b) is shown in Fig. 7 (c). Connected component analysis is performed over B_m to extract the regions with high probability to belong to vehicles. Resulting regions are filtered according to a minimum area criterion in order to remove noise.

Naturally, regions corresponding to the tracked vehicles should exist in B_m . Besides, if there is some additional region in B_m this is regarded as a potential new vehicle in the image and it is further analyzed in the hypothesis verification stage. In particular, in the example in Fig. 7 (c) three regions are obtained: the upper two regions correspond to existing vehicles, labeled 1 and 2, while the small region in the lower left corner

constitutes a potential new vehicle (in this case it is actually a vehicle, as can be observed in Fig. 7 (a)). Since only the lower part of the vehicles is reliable in the rectified domain, candidates are characterized by their position and width. As potential vehicles are verified according to their appearance in the original image, their position and width in this domain are computed by means of the inverse transformation T^{-1} . Finally, a 1:1 aspect ratio is initially assumed for the vehicle so that a bounding box R_h can be hypothesized for vehicle verification.

7.2 Hypothesis verification

Vehicle verification is based on a supervised classification stage based on Support Vector Machines (SVM). A database of vehicle rear images is generated for the training of the classifier as will be explained in Section 7.2.2. Most importantly the database separates images according to the region in which the vehicle is found (close/middle range in the front, close/middle range in the left, close/middle range in the right, and far range). Indeed, the view of the vehicle rear changes in these areas and thus affects its intrinsic features. This is taken into account in the design of the feature description, which adapts to the particularities of the different areas. Besides, a different classifier is trained for each of them using the corresponding subsets of images in the database.

As for the feature description, a new descriptor is proposed based on two of characteristics that are inherent to the vehicles: high edge content

and symmetry. Indeed, the method automatically adapts the area for feature extraction according to a vertical symmetry-based local window refinement. This allows to correct position offsets in the hypothesis generation stage and to adapt to the vehicle rear contour. Regarding the feature extraction within the refined region, a new descriptor that exploits the inherently rectangular structures of the vehicle rear is designed. The descriptor, called CR-HOG, is based on the analysis of Histograms of Oriented Gradients (HOGs) in concentric rectangles around the center of symmetry, called CR-HOG.

7.2.1 CR-HOG feature extraction: HOGs evaluate local histograms of image gradient orientations in a dense grid. The underlying idea is that the local appearance and shape of the objects can often be well characterized by the distribution of the local edge directions, even if the corresponding edge positions are not accurately known.

This idea is implemented by dividing the image into small regions called cells. Then, for each cell, a histogram of the gradient orientations over the pixels is extracted. The original HOG technique, proposed by Dalal and Triggs [12], presents two different kinds of configurations, called Radial HOG (R-HOG) and circular HOG (C-HOG), depending on the geometry of the used cells. Specifically, the former involves a grid of rectangular spatial cells and the latter uses cells partitioned in a log-polar fashion.

As stated, in this work we present a new configuration for the cells that better adapts the characteristics of the vehicles. Indeed, the rear of the vehicles presents an inherently rectangular structure: not only is the

outer contour of the vehicle rear quasi-rectangular, but the inner structures such as the license plate and the rear window are also rectangular. Hence, we naturally define a new configuration of HOG composed of concentric rectangular cells as shown in Fig. 8 (a). This structure will be referred to as CR-HOG (for concentric rectangle-based HOG). The layout of the CR-HOG has five parameters: the number of concentric rectangles n , the number of orientation bins b , the center c_s of the window, its height h_s , and its width w_s .

In practice, the hypothesized region for vehicle verification, R_h , may not perfectly fit the actual bounding box of the vehicle in terms of size and alignment. In particular, it is often the case that the vehicle is not perfectly centered in R_h , especially in the horizontal axis. Therefore, direct application of CR-HOG (or of standard HOG) over R_h will possibly result in degraded performance. Instead, we refine the region likely containing the vehicle through the analysis of vertical symmetry in the intensity of the region. In particular, the subregion within H_s giving the maximum degree of vertical symmetry is kept for HOG computing. Vertical symmetry is calculated using the method in [40]. As a result, we obtain the axis of vertical symmetry, x_s and the width of the region that maximizes the symmetry measure, w_s . The height h_s of the window for HOG application is taken as that of R_h and its center is thus given by $c_s = (x_s, h_s/2)$.

Fig. 8 (b) illustrates the window adaptation approach based on symmetry analysis. Observe that the refined vertical side limits (painted in red)

fit much better the bounding edges of the vehicle rear. In practice, the area for calculation of CR-HOG is extended by a 10% so that the outer edges of the vehicle are also accounted for in the descriptor.

The steps for the calculation of CR-HOG on the refined window are the following. First, the gradient magnitude and orientation is computed at each point of the window using some standard operator (Sobel 3×3 masks are used in our implementation). Then, in order to create a histogram of orientations, a number of orientation bins is defined and each pixel votes for the bin that contains its corresponding angle. The votes are weighted by the magnitude of the gradient at that point. Three possible configurations have been considered involving $n = 8, 12$ and 18 bins evenly spaced over $[0, 180)$, as illustrated in Fig. 9.

The bins have been shifted so that the vertical and horizontal orientations, which are very frequent in the rear of vehicles, are in the middle of their respective bins. This way, small fluctuations around 0° and 90° will not affect the descriptor. The histogram of each cell is finally normalized by the area of the cell so that histograms of different cells are in the same order of magnitude. The CR-HOG descriptor is composed of the normalized histogram of orientations of all the cells. The optimal configuration of the number of orientation bins, n , and the number of cells, b , is discussed in Section 8 for each region of the image.

7.2.2 Classification stage: The CR-HOG descriptors are introduced in a Support Vector Machine-based classifier. A new database containing 4000

positive vehicle images and 4000 negative vehicle images is used to train and test the classifier (this is accessible at <http://www.gti.ssr.upm.es/~jal>). The core of the database is composed of images of our own collection; besides, images have also been extracted from the Caltech [41,42] and the TU Graz-02 [43,44] databases and included in the data set. The joint database consists of images of resolution 64×64 acquired from a vehicle-mounted forward-looking camera. Each image provides a view of the rear of a single vehicle. Some images contain the vehicle completely while others have been drawn to contain it only partially (all images contain at least 50% of the vehicle rear) in order to simulate putative results of the hypothesis generator.

Images involve many different viewpoints of the vehicle rear corresponding to vehicles in different locations relative to the vehicle in which the camera is mounted. Specifically, the space is divided in four main regions: close/middle range in the front, close/middle range in the front, close/middle range in the right, and far range. The database contains 1000 images of each of these views. A set of 4000 images not containing vehicles have also been used to train and test the classifier. The images are selected in such a way that the variability in terms of vehicle appearance, pose, and acquisition conditions (e.g., weather conditions, lighting) is maximized. A classifier based on SVM using linear basis functions is used for each of the four image regions.

8 Experiments and discussion

Experiments regarding the proposed method have been performed two-fold. On the one hand, the performance of the novel CR-HOG based approach for vehicle detection is tested on the database referred to in Section 7.2.2. On the other hand, experiments have been carried out in the complete system integrating vehicle detection and tracking over a wide variety of video sequences.

The SVM-based classifier for vehicle detection explained in Section 7 has been trained and tested in Matlab using the Bioinformatics Toolbox. The method involves two design parameters, namely the number of orientation bins in the histogram, n , and the number of cells, b . Experiments have been performed on the database for values $n = 8, 12, 18$ and $b = 2, 3, 4$. Cross-validation procedure is used to test the method. Specifically, 50% of the images are randomly selected for the training set and the remaining 50% is used as the testing set. This process is repeated 5 times and the average is computed.

The accuracy or correctly classified rate of samples as a function of these parameters is provided in Table 1 for each of the four regions. These results are graphically represented in Fig. 10 and Fig. 11 to facilitate their interpretation. In particular, Fig. 10 shows the accuracy results as a function of the number of orientation bins, n , by averaging on b , for each area of the image. Analogously, Fig. 11 illustrates the average accuracy results as a function of the number of cells, b .

As a first conclusion of the experiments we infer that the accuracy decreases for $b = 4$ in all the areas of the image. As for the number of orientation bins, a different behavior is observed for the frontal and the sides views. Specifically, for the central close/middle and far ranges similar results are obtained for $n = 8$ and $n = 12$, while the performance decreases notably for $n = 18$. As opposed to it, for the left and right areas a significative accuracy increase is observed from $n = 8$ to $n = 12$; a further increase to $n = 18$ does not bring an additional gain. This contrast is indeed reasonable, since from a completely orthogonal viewpoint the edges of a vehicle are fairly invariant and mostly vertical and horizontal; conversely, in the side views the upright edges corresponding to the back window and its contour (especially the furthest from the image center) tend to divert from verticality. Consequently more variability is found in the gradient orientation map, and therefore more bins are necessary to capture fine-detail.

A good trade-off between complexity and performance is achieved by selecting $(b, n) = (2, 8)$ for the close/middle and far ranges, and $(b, n) = (3, 12)$ for the left and right views. This involves respective detection accuracies of 94.88%, 85.92% 91.82%, and 89.42%, which results in an average correct detection rate of 90.51%. The rate difference between left and right views is due to the particularities of the traffic participants in the right lane, which usually includes slow vehicles (buses, trucks, vans, etc.), which involve a great appearance variety and hence make classification much more challenging. Naturally the worse classification rate is obtained for the further

vehicles, in which the edge-details are degraded. The results are improved to an overall accuracy of 92,77% when using a Gaussian radial basis function kernel (instead of the linear kernel), with respective correct detection rates of 96.14%, 89.92%, 94.14% and 90.86% for the different areas. However, the proposed method continuously generates hypotheses for the potential vehicles, hence, even if a vehicle is not detected in a given frame, it is usually detected in the following frames. Therefore, the small latency incurred by the linear kernel-based classification is usually negligible and it is not necessary to use more complex kernels.

As regards vehicle tracking, the designed method has been tested on a wide variety of sequences recorded on Madrid, Brussels and Turin. These sequences, which were acquired in several sessions with different driving conditions (i.e., illumination conditions, weather, pavement color, traffic density, etc.) amount to 22 : 38 minutes. Test sequences were acquired from a forward-looking digital video camera installed near the rear mirror of a vehicle driven in highways. The method is able to operate near real-time at 10 fps on average over an Intel(R) Core(TM) i5 processor running at 2,67 GHz. Implementation is done in C++.

The above-mentioned test sequences are used to compare the performance of the proposed tracking method with the two methods most widely used in the literature. Namely, those involve independent tracking of multiple objects with a Kalman filter assigned to each object (which will be shortly referred to as KF-based tracking), or joint tracking using parti-

cles filters based on importance sampling (shortly, SIS-based tracking). For the implementation of KF-based tracking, appearance-based region labeling through connected-component analysis is used as in Section 7.1 to locate vehicles in every frame, and tracks are formed temporally by matching the regions according to a minimum distance criterion. As for SIS-tracking, a sequential resampling scheme is used (see details of the algorithm in [21]). Additionally, the motion model used for SIS-based tracking is exactly that designed for the proposed method, while KF-based tracking uses the same dynamic model for independent motion of vehicles, but cannot accommodate any interaction model. Other parameters of the dynamic model are $\sigma_q = (10, 15)$, $w_l = 90$ and $d_s = 96$. Regarding the observation model, the dimensions of the local windows R_a and R_b are set to $w = h = 10$. Also, the standard deviation of the proposal density is optimally calculated for the proposed method and for SIS-based tracking as $\sigma_q = (2, 3)$ and $\sigma_q = (5, 8)$, respectively. Finally, the same number of samples $N = 250$ is used for both methods.

To compare methods, the number of tracking failures incurred by each of the methods on the same test sequences is counted. A tracking failure is considered when the tracker fails to provide continuous and coherent measures for a given vehicle inside the Region of Interest (ROI). The ROI is defined to be the scope of the Inverse Perspective Mapping, which usually comprises the own and the two adjacent lanes, and extends longitudinally up to a distance d_f that depends on the camera calibration. Comparative

results are displayed in Table 2. As expected, the proposed method largely outperforms the others in terms of tracking failures in the test sequences. Naturally, KF-based tracking delivers the highest error rate as it is unable to deal with situations in which vehicles interact. Notably, SIS-based tracking also incurs in a significant number of errors, since the number of particles is relatively small and fail to correctly sample the space when the number of vehicles grows.

The strength of the method lies to a great extent in the combination of two different sources of information (appearance and motion) for the definition of the observation model. Indeed, the combination of information ensures that whenever the two sources are available a robust average estimate is produced, and most importantly, it allows to keep track of the objects even if one of the information sources is unavailable or unreliable. Fig. 12 (a) illustrates the sampling process for the original image in Fig. 12 (a1) whenever the two types of information are available. In particular, Fig. 12 (a2) shows the rectified domain after IPM application, Fig. 12 (a3) corresponds to the appearance-based pixel-wise classification (in which pixels likely belonging to the lower parts of vehicles are painted in white as explained in Section 5.2), and Fig. 12 (a4) contains analogously the pixel-wise motion-based classification as explained in Section 6.2.

The process of generation of samples in the framework of the Markov chain is superimposed on Fig. 12 (a3) and Fig. 12 (a4). In particular, the segment between the previous sample and the proposed sample is painted in

green whenever the latter is accepted, and in red if it is rejected. As can be observed, accepted samples concentrate in the area of high likelihood (i.e. the transition between road and vehicles), while samples diverging from this area are rejected. The final estimates for vehicles positions are indicated in Fig. 12 (a5) with white segments underlining the vehicle rears.

As stated, dual modeling from two sources prevents the method from losing track whenever one of the sources is unreliable. This is illustrated in Fig. 12 (b) and Fig. 12 (c), where the sampling process is depicted analogously to Fig. 12 (a) for the images in Fig. 12 (b1) and Fig. 12 (c1). In particular, in Fig. 12 (b) the motion-based observation provides no measurement for the right vehicle (Fig. 12 (d)), however this is compensated by the correct appearance-based observation, which avoids particle dispersion. Therefore, the vehicle is correctly tracked as shown in Fig. 12 (b5). In contrast, the particles for the left vehicle overcome an inaccurate initialization and converge to its actual position due to good appearance-based and motion-based observations. The opposite case is illustrated in Fig. 12 (c), in which the appearance-based model fails to detect the furthest vehicle (see Fig. 12 (c3)), whereas the region of motion is still observable in the difference between aligned images in Fig. 12 (c4). This allows to keep track of the vehicle, as shown in Fig. 12 (c5).

Finally, apart from the statistical results in Table 2, some graphical examples of the performance of the method are shown in Fig. 13. This figure displays snapshots of the tracking process for four different time points

(from left to right), for three different example sequences. In the first sequence the method simultaneously tracks a vehicle that is being rapidly overtaken by the own vehicle. Most interestingly, there is some degree of interaction between the vehicles, in fact, they are fairly close in Fig. 13 (a3). In traditional PF-based tracking methods particles are prone to concentrate around the object with the highest likelihood which results in the loss of the other object. In contrast, the designed interaction model allows to prevent this situation and to successfully track the vehicles until they part. In the second example, in Fig. 13 (b), tracking of a vehicle driving at a slow pace in the right lane is shown. At the same time, the method swiftly detects a vehicle in the left hand side (Fig. 13 (b3)) and tracks it until it is far away, while also keeping track of the vehicle in front of the own car. Finally, in Fig. 13 (c) simultaneous tracking of several vehicles is shown: the vehicle ahead the own car moves from the lower-left corner of the image to the upper-middle part, while at the same time tracking is kept for the distant vehicle in the right lane. Meanwhile, a new vehicle entering the scene in the left hand is detected and tracked until it is nearly at the same distance as the other vehicles.

9 Conclusions

In this paper a new probabilistic framework for vehicle detection and tracking has been presented based on MCMC. As regards vehicle detection, a new descriptor, CR-HOG, has been defined based on the extraction of gra-

dient features in radial rectangular bins. The descriptor has proven to good discriminative properties using a reduced number of features in a simple linear-kernel SVM classifier, and is thus ideally suited for real-time applications. In addition, the tracker is proven to perform better than state of the art methods based on Kalman and particle filtering in terms of tracking failures. The power of the algorithm lies on the fusion of information of different nature, especially regarding the observation model. In effect, apart from appearance the method exploits another feature that is inherent to vehicles, their motion, through the analysis of the geometry between successive views of the scene. In addition, MCMC method is exploited to perform efficient sampling and to avoid the performance degradation of particle filter-based approaches in multiple object tracking arising from the curse of dimensionality. In summary, the method is able to overcome the usual limitations of particle filter-based approaches and to provide robust vehicle tracking in a wide variety of driving situations and environment conditions.

10 Acknowledgements

This work was supported in part by the Ministerio de Ciencia e Innovación of the Spanish Government under projects TEC2007-67764 (SmartVision) and TEC2010-20412 (Enhanced 3DTV).

References

1. GY Song, KY Lee, JW Lee, in Proceedings of the IEEE Intelligent Vehicles Symposium, Eindhoven, The Netherlands, 4–6 June 2008.
2. L Gao, C Li, T Fang, Z Xiong, in Image Analysis and Recognition, ed. by A Campilho, M Kamel. International Conference on Image Analysis and Recognition, Pvoa de Varzim, June 2008. Lecture Notes in Computer Science, vol. 5112 (Springer, Heidelberg, 2008), 142–150.
3. L-W Tsai, J-W Hsieh, K-C Fan, IEEE Transactions on Image Processing **16**(3), 850–864 (2007).
4. W Liu, X Wen, B Duan, H Yuan, N Wang, in Proceedings of the IEEE Intelligent Vehicles Symposium, Istanbul, Turkey, 13–15 June 2007.
5. A Ess, B. Leibe, K Schindler, L van Gool, IEEE Transactions on Pattern Analysis and Machine Intelligence **31**(10), 1831–1846 (2009).
6. G Toulminet, M Bertozzi, S Mousset, A Benschraier, A Broggi, IEEE Transactions on Image Processing **15**(8), 2364–2375 (2006).
7. TN Tan, IEEE Transactions on Image Processing, **9**(8), 1343–1356 (2000).
8. JM Collado, C Hilario, A de la Escalera, JM Armingol, Proceedings of the IEEE Intelligent Vehicles Symposium, University of Parma, Italy, 17 June 2004.
9. DM Ha, J-M Lee, Y-D Kim, Image and Vision Computing **22**, 899–907 (2004).
10. Z Sun, G Bebis, R Miller, in Proceedings of the International Conference on Digital Signal Processing, Santorini, Greece, 1–3 July 2002.
11. Z Sun, G Bebis, R Miller, IEEE Transactions on Image Processing **15**(7), 2019–2034 (2006).

12. N Dalal, B Triggs, in Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, San Diego, California, 20–26 June 2005.
13. P Negri, X Clady, SM Hanif, L Prevost, A cascade of boosted generative and discriminative classifiers for vehicle detection. *EURASIP Journal on Advances in Signal Processing* (2008). doi:10.1155/2008/782432.
14. C-CR Wang, J-JJ Lien, *IEEE Transactions on Intelligent Transportation Systems*, **9**(1), 83–96 (2008).
15. C Papageorgiou, T Poggio, *International Journal of Computer Vision*, 38(1), 15–33 (2000).
16. R Lienhart, J Maydt, in Proceedings of the IEEE International Conference on Image Processing, Rochester, New York, 22–25 Sept. 2002.
17. J Arrospe, L Salgado, M Nieto, F Jaureguizar, in Proceedings of the IEEE International Conference on Image Processing, San Diego, California, 12–15 Oct. 2008.
18. Y Chen, M Das, D Bajpai, in Proceedings of the Canadian Conference on Computer and Robot Vision, Montreal, Quebec, 28–30 May 2007.
19. R Goecke, N Pettersson, L Petersson, in Proceedings of the IEEE Intelligent Vehicles Symposium, Istanbul, Turkey, 13–15 June 2007.
20. M Asadi, F Monti, CS Regazzoni, Feature classification for robust shape-based collaborative tracking and model updating. *EURASIP Journal on Image and Video Processing* (2008). doi:10.1155/2008/274349.
21. MS Arulampalam, S Maskell, N Gordon, T Clapp, *IEEE Transactions on Signal Processing*, **50**(2), 174–188 (2002).
22. A Dore, M Soto, CS Regazzoni, *IEEE Signal Processing Magazine* **27**(5), 46–55 (2010).

23. JJ Pantrigo, A Sánchez, AS Montemayor, A Duarte, *Pattern Recognition Letters* **29**, 1160-1174 (2008).
24. T Gao, Z-G Liu, W-C Gao, J. Zhang, in *Advances in Neuro-Information Processing*, ed. by M Kppen, N Kasabov, G Coghill. 15th International Conference on Neural Information Processing of the Asia-Pacific Neural Network Assembly, Auckland, Nov. 2008. *Lecture Notes in Computer Science*, vol. 5507 (Springer Berlin, Heidelberg, 2008), 695–702.
25. C Luo, X Cai, J Zhang, in *Proceedings of IEEE Workshop on Multimedia Signal Processing*, Cairns, Australia, 8–10 Oct. 2008.
26. J Wang, Y Ma, C Li, H Wang, J Liu, in *Proceedings of World Congress on Computer Science and Information Engineering*, Los Angeles, California, 31 March–2 April 2009.
27. Z Khan, T Balch, F Dellaert, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **27**(11), 1805–1819 (2005).
28. J Wang, Y Yin, H Man, Multiple human tracking using particle filter with Gaussian process dynamical model. *EURASIP Journal on Image and Video Processing* (2008). doi:10.1155/2008/969456.
29. K Smith, D Gatica-Perez, J-M Odobez, in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, San Diego, California, 20–26 June 2005.
30. CM Bishop, *Pattern Recognition and Machine Learning* (Springer, New York, 2006).
31. M Bertozzi, A Broggi, *Computer Vision and Image Understanding* **113**(6), 743–749 (1998).
32. R Danescu, S Nedeveschi, M-M Meinecke, and TB To, in *Proceedings of the IEEE Intelligent Vehicles Symposium*, Eindhoven, The Netherlands, 4–6 June

- 2008.
33. M Nieto, L Salgado, F Jaureguizar, J Cabrera, in Proceedings of the IEEE Intelligent Vehicles Symposium, Istanbul, Turkey, 13–15 June 2007.
 34. JA Bilmes, A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models, Technical Report TR-97-021 (1998).
 35. RI Hartley, A Zisserman, *Multiple View Geometry in Computer Vision* (Cambridge University Press, Cambridge, 2000).
 36. C Harris, M Stephens, in Proceedings of the 4th Alvey Vision Conference, Manchester, England, 31 Aug.–2 Sept. 1988.
 37. BD Lucas, T Kanade, in Proceedings of the 7th Joint Conference on Artificial Vision, Vancouver, Canada, 24–28 Aug. 1981.
 38. DG Lowe, International Journal of Computer Vision **60**(2), 91–110 (2004).
 39. TK Moon, WC Stirling, *Mathematical Methods and Algorithms for Signal Processing* (Prentice Hall, Englewood Cliffs, NJ, 1999).
 40. T Zielke, T Brauckmann, W Von Seelen, CVGIP: Image Understanding, **58**(2), 177–190 (1993).
 41. The Caltech Database (Computational Vision at California Institute of Technology, Pasadena), <http://www.vision.caltech.edu/html-files/archive.html>. Accessed 14 May 2011.
 42. R Fergus, P Perona, A Zisserman, in Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Madison, Wisconsin, 16–22 June 2003.
 43. The TU Graz-02 Database (Graz University of Technology), <http://www.emt.tugraz.at/~pinz/data/GRAZ.02/>. Accessed 14 May 2011.

44. A Opelt, A Pinz, in Proceedings of the 14th Scandinavian Conference on Image Analysis, Joensuu, Finland, 19–22 June 2005.

Fig. 1 General scheme of the proposed method.

Fig. 2 Transformation to the rectified domain through Inverse Perspective Mapping. As opposed to the original image (a), in the rectified image (b) the effect of perspective is removed and thus motion of vehicles is easier to model.

Fig. 3 Steps for initialization of EM regarding intensity feature parameters. The sequence of images is the following: (a) example image in rectified domain; (b) binary mask for high-gradient pixel removal (pixels in white are removed); (c), (d) and (e) maps of pixels for pavement, lane marking, and vehicle class characterization, respectively.

Fig. 4 Histograms of the (a) pavement, (b) lane marking, and (c) vehicle class pixel maps in Fig. 3.

Fig. 5 Relative pose of the camera at two different time points k_1 and k_2 . The world coordinate system has its origin at the position of the camera center at k_1 .

Fig. 6 Example of image alignment. Image (a) and (b) correspond to times $k-1$ and k , respectively, of the video sequence; (c) is the image at $k-1$ warped with \hat{H}^k ; (d) is the difference between aligned images, i.e., (b) and (c); and (e) is the corresponding image of difference in the rectified domain. Both (d) and (e) have been binarized for better visualization: white regions correspond to regions of difference, which usually correspond to the lower edge of vehicles.

Fig. 7 Example of generation of a new vehicle hypothesis. The sequence of images is the following: (a) original image, (b) rectified image, (c) binary map B_m corresponding to appearance analysis in (b): pixels in white indicate potential location of vehicles. In the example, the regions labeled 1 and 2 correspond to existing vehicles, while the small region arising in the lower left corner constitutes a potential new vehicle.

Fig. 8 Combined HOG and symmetry based descriptor. In (a) the structure of concentric rectangle HOG (CR-HOG) with its corresponding parameters is illustrated. In (b) the refined region obtained after vertical symmetry analysis is shown: green and red lines indicate respectively the symmetry axis and the width of the region yielding the maximum symmetry values.

Fig. 9 Possible configurations of CR-HOG regarding the number of orientation bins. The range of gradient orientation angles [0-180) is divided in uniformly spaced sectors. Pixels with gradient orientations inside each sector accumulate to the corresponding bin of the histogram proportionally to the magnitude of their gradient. Configurations with (a) 8, (b) 12, (c) 18 bins are considered.

Fig. 10 Classification accuracy as a function of the number of cells, b . The results are broken down for images corresponding to (a) close/middle, (b) left, (c) right and (d) far views and for $b = 2, 3$ and 4.

Fig. 11 Classification accuracy as a function of the number of orientation bins, n . The results are broken down for images corresponding to (a) close/middle, (b) left, (c) right and (d) far views and for $n = 8, 12$ and 18.

Fig. 12 Illustration of sampling process for different example images. From left to right, images correspond to the (1) original image, (2) rectified domain, (3) appearance-based vehicle probability map, B_m , (4) motion-based vehicle probability map, and (5) tracking results. The sampling process is illustrated in images (3) and (4): accepted and rejected particles are painted in green and red, respectively. Images (2)-(4) are zoomed for better visualization of the sampling process. Images in (a) illustrate a normal sampling scenario, while images in (b) and (c) show how combined sampling is able to overcome bad (b) motion-based and (c) appearance-based measurements.

Fig. 13 Vehicle tracking for three different sequences (a)-(c). From left to right, the images show results at times $k_0, k_0 + 200, k_0 + 340, k_0 + 440$; $k_0, k_0 + 170, k_0 + 215, k_0 + 295$; $k_0, k_0 + 250, k_0 + 360, k_0, k_0 + 460$ for sequences (a), (b), and (c), respectively

Table 1 Classification accuracy rates of CR-HOG

	Middle/Close			Left			Right			Far		
	$\beta = 2$	$\beta = 3$	$\beta = 4$	$\beta = 2$	$\beta = 3$	$\beta = 4$	$\beta = 2$	$\beta = 3$	$\beta = 4$	$\beta = 2$	$\beta = 3$	$\beta = 4$
$n = 8$	94,88	94,98	94,68	91,04	91,18	91,16	88,58	89,14	87,94	85,92	85,86	85,76
$n = 12$	94,96	94,80	95,14	91,46	91,82	91,46	89,28	89,42	88,16	85,32	85,24	85,16
$n = 18$	94,78	93,96	93,24	91,98	91,60	91,06	89,34	88,84	88,10	85,76	85,22	84,60

Table 2 Summary of tracking results

Method	Tracking failures	Number of frames	Number of vehicles
KF-based Tracking	36		
PF-based Tracking	31	33454	120
Proposed Method	9		

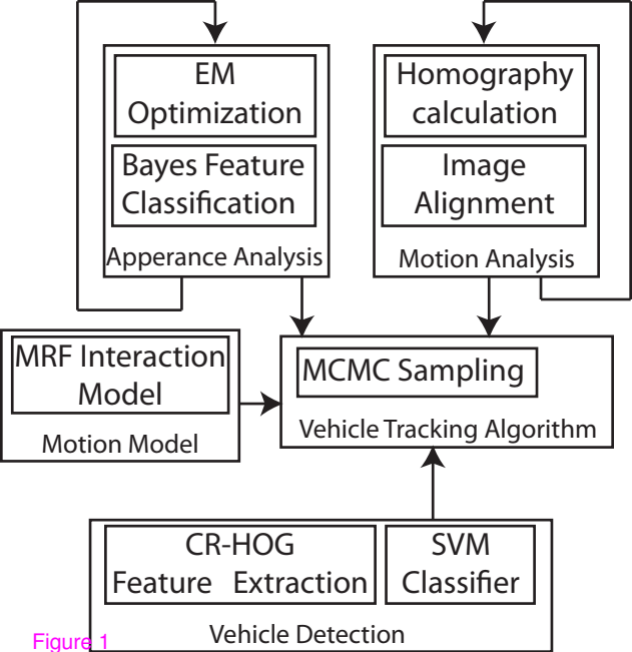
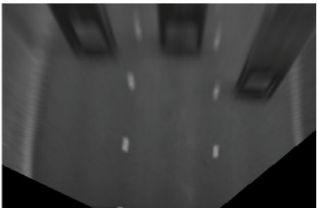


Figure 1



Figure 2 (a)



(b)

Rectified plane



Edge mask



Pavement pixels



Lane markings pixels



Object pixels



Figure 3 (a)

(b)

(c)

(d)

(e)

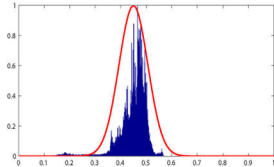
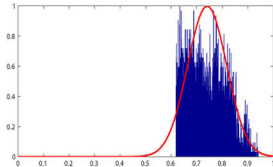
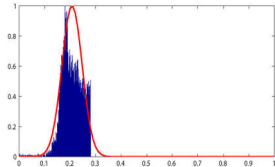


Figure 4 (a)



(b)



(c)

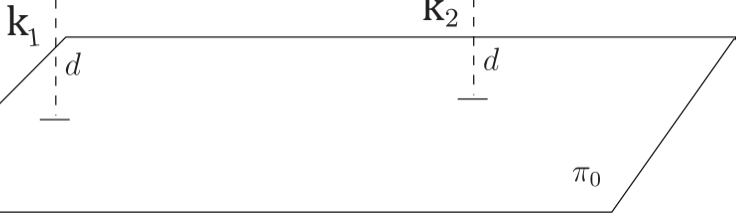
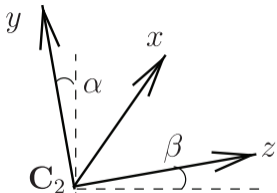
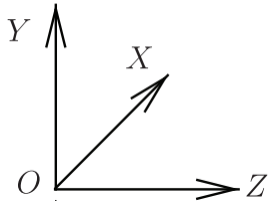


Figure 5



(a)



(b)



(c)



(d)

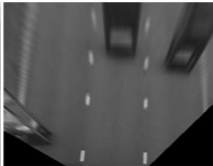


(e)

Figure 6



Figure 7(a)



(b)



(c)

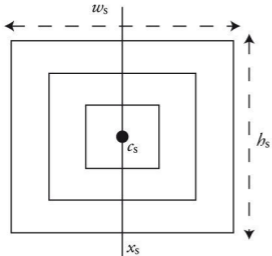
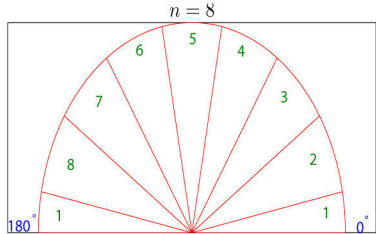
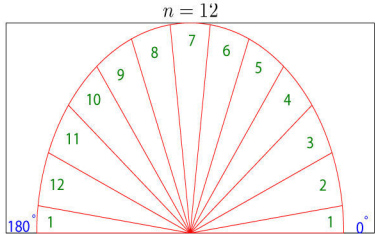


Figure 8 (a)

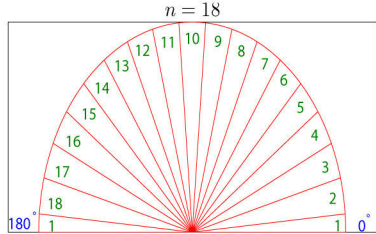
(b)



(a)

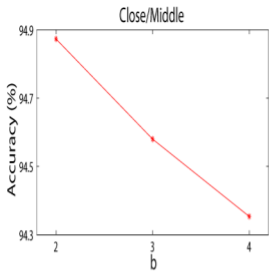


(b)

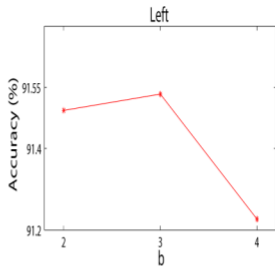


(c)

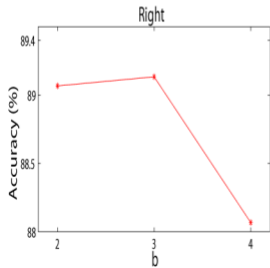
Figure 9



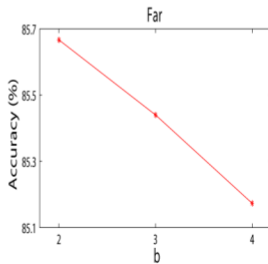
(a)



(b)

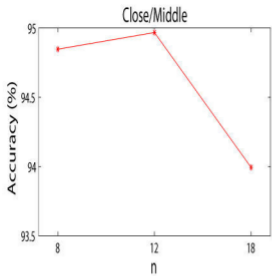


(c)

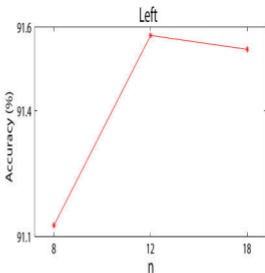


(d)

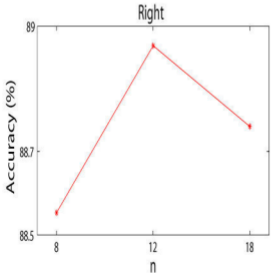
Figure 10



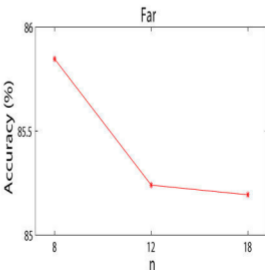
(a)



(b)



(c)



(d)

Figure 11

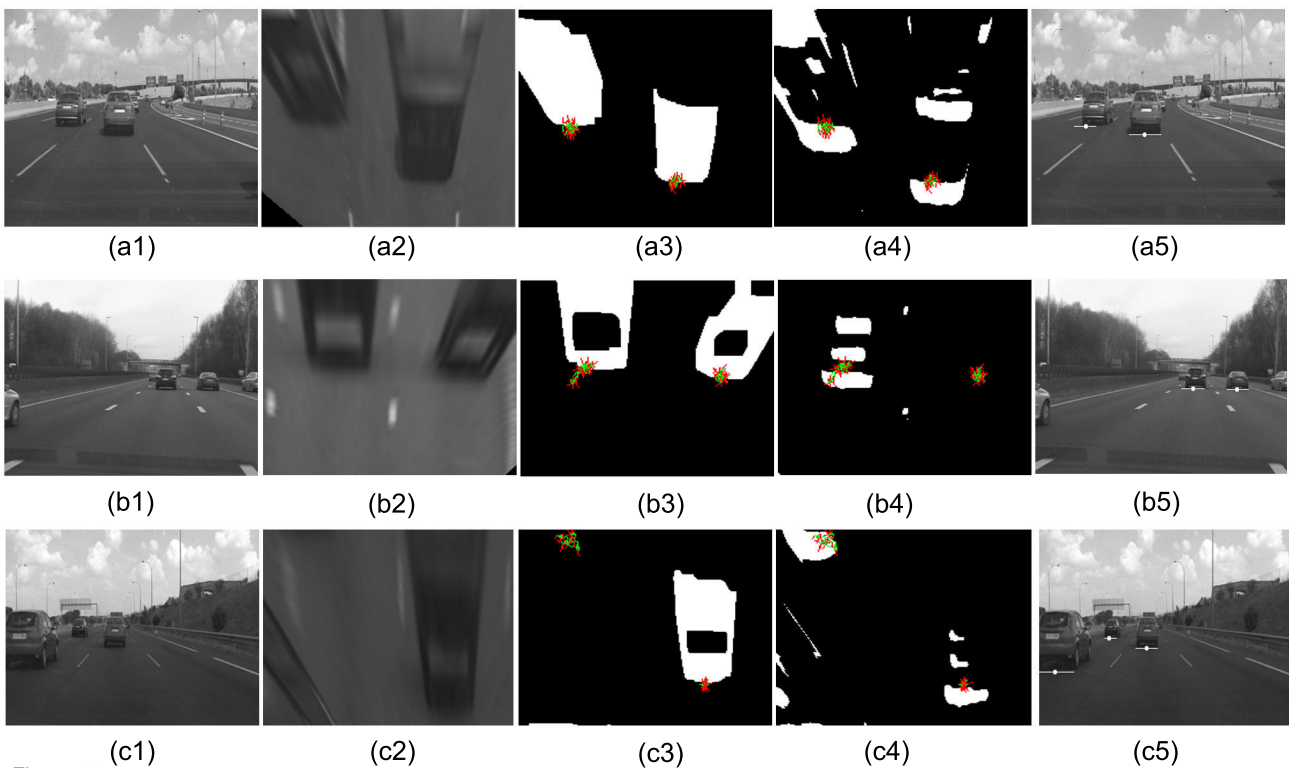


Figure 12



(a1)



(a2)



(a3)



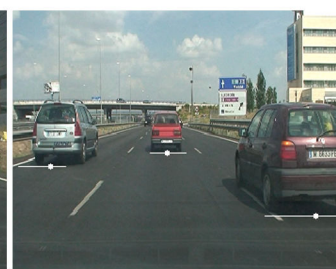
(a4)



(b1)



(b2)



(b3)



(b4)



(c1)



(c2)



(c3)



(c4)

Figure 13