

Sensory Augmentation for Increased Awareness of Driving Environment

John M. Dolan

Paul Rybski

Dec. 14, 2013

Technologies for Safe and Efficient Transportation (T-SET) UTC

The Robotics Institute

Carnegie Mellon University

Pittsburgh, Pennsylvania 15213

DISCLAIMER

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. This document is disseminated under the sponsorship of the U.S. Department of Transportation's University Transportation Centers Program, in the interest of information exchange. The U.S. Government assumes no liability for the contents or use thereof.

Table of Contents

Problem	4
Approach/Methodology	4
Findings	5
Conclusions/Recommendations	9
Works Cited	10
IV 2012 Paper (Pedestrian Detection).....	11

Problem

The goals of this project were to extend the state of the art of vehicle perception systems for use in roadway traffic and develop systems that can model and predict the actions of multiple simultaneous road users so as to identify potentially hazardous situations before they can turn into accidents.

Approach/Methodology

We augmented vehicles with sensors and processing capabilities to perceive obstacles (both static and dynamic), predict how those obstacles might move over time, identify locations where unseen hazards might appear, and continually evaluate these values to determine the possibility that an unsafe condition might occur in the immediate future. While the Urban Challenge (Urmson & et al., 2008) focused on fully autonomous vehicles, similar perception systems can also be deployed in manually-driven cars that could alert the human driver if an unsafe road condition is approaching. We used behavioral models of traffic to identify the perceived intent of nearby vehicles, use those intent models to predict the most likely future positions of those vehicles, and determine whether a potentially unsafe condition may arise in the near future.

For a vehicle to automatically predict unsafe situations that may occur in traffic, it needs to rely on successfully detecting, tracking, *modeling*, and *predicting* the motions of other moving objects (e.g. cars, bicyclists, and pedestrians) in its surroundings. We therefore developed novel and robust approaches for the modeling and prediction of vehicle motion. Once each object's intent has been identified, these behavior models can be used to predict a series of "idealized" potential future trajectories. Each of these future trajectories is weighted by the likelihood of its occurring as well as by the potential to cause an unsafe traffic condition. With this information, our perception system can continuously measure the risk of the current situation's turning into an

unsafe situation that could end up causing an accident. Similarly, by reasoning about the known roadmap on which the vehicles and other road users are operating, evaluating the visible and blind areas of the sensors (e.g. blind corners and obstacles blocking sensors' views), the perception system can alert the driver that he is approaching a potentially unsafe situation in time for him to take appropriate actions to remove himself from the situation (e.g. slow down, change lanes, pull over, etc.).

Findings

The period of performance for this project was February 2012 – December 2013. Paul Rybski was the PI during 2012, and upon his taking a leave of absence to work at Caterpillar, John M. Dolan took over as PI during 2013. The report's findings are therefore in two major categories: pedestrian detection work performed during 2012 under Dr. Rybski's supervision and traffic detection work performed during 2013 under Dr. Dolan's supervision.

Pedestrian detection. The goal of this work was the creation of a real-time pedestrian detection system suitable for use on automotive-grade camera and computing platforms. The developed system demonstrated superior performance when compared to many state-of-the-art detectors and was able to run at 14 fps on an Intel Core i7 computer when applied to 640x480 images. The system demonstrated a 61% detection rate on the largest publicly available pedestrian detection dataset, the Caltech Pedestrian Benchmark, whereas recent other state-of-the-art detectors achieve a generally lower detection rate between 50% and ~61%, but at significantly greater computational expense. These findings resulted in a paper presented at the 2012 IEEE Intelligent Vehicles Symposium (Cho, Rybski, & Zhang, 2012), and also included at the end of the current report.

Traffic detection. In this work, our goal was to develop algorithms for traffic detection and response capable of using relatively inexpensive and low-capability automotive-grade sensors compared to the high-cost, high-capability sensors used more typically on robotic vehicles. A good example of the latter is the \$70,000 Velodyne spinning laser sensor that most of the 2007 Urban Challenge vehicles used in order to get a high-density 3D range point cloud. Such a sensor is not practical for production automobiles from a cost or appearance/integration standpoint. Our lab's Cadillac SRX vehicle instead uses IBEO LiDAR and automotive-grade radar for traffic detection. The IBEO sensors are not available for automobiles, but they are much lower-cost than the Velodyne (several thousand dollars vs. \$70,000), and a very similar model is planned for release in an automotive-grade form by the French automotive components manufacturer Valeo (<http://www.valeo.com>) within the next two years.

The typical approach to sensor fusion on an autonomous or semi-autonomous vehicle is to implement tracking, which attempts to correlate data from cycle to cycle to obtain a picture of all “targets” in sensor range and their state over time. This can be computationally intensive and typically involves lags which are undesirable when attempting to make rapid decisions in traffic. Our approach in this project has been instead to use computationally lightweight heuristics aimed at determining the essential information for particular traffic situations, specifically lane changing and merging.

Lane changing. Figure 1 shows the relevant considerations for our lane change heuristic. As the caption describes, there are three relevant considerations. We first use LiDAR to check whether a polygon sufficiently large to permit lane change is free in the intended lane. LiDAR is more suited for this than radar because of its better persistence (i.e., lesser intermittency) when obstacles are present. We then use radar to check the velocity of the closest following vehicle in

the intended lane in order to determine whether the time to collision is sufficiently greater than the time to change lanes. If the time to collision is infinite, the SRX and the following are traveling at the same speed, and there is no problem. If the following vehicle's speed is greater than the SRX's, there is a finite time to collision, and we ensure that it is long enough for the following vehicle to recognize the SRX's intention and reduce speed accordingly. Finally, we use radar to determine based on the relative speeds of the SRX and any leading vehicle whether distance-keeping (adaptive cruise control) can be safely invoked as and after the lane change is performed.

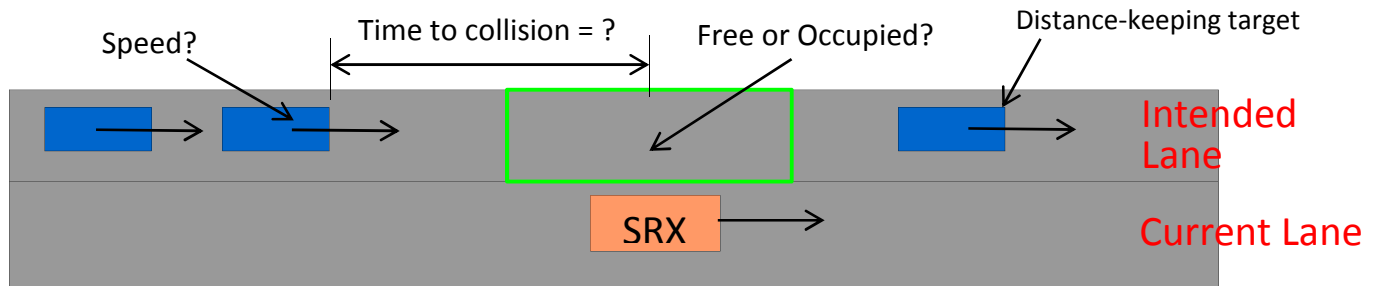


Figure 1. Lane change considerations for case-specific traffic detection algorithm. SRX is the ego vehicle, which is checking whether a move from the Current to the Intended Lane is feasible. The three basic considerations are: a) Is the slot the ego vehicle wants to move into free or occupied? b) Is the time to collision for the closest following vehicle sufficiently greater than the time to change lanes? c) Is the leading vehicle far enough ahead for safe distance keeping upon making the lane change?

Note that this method uses the particular strengths of each sensor modality (LiDAR and radar) without explicitly fusing them into a computationally more complex and expensive tracking system. Figure 2 shows an example from actual vehicle data of the use of the two types of sensor: on the left, the LiDAR sensors determine there is a vehicle in the polygon, so a lane change is currently infeasible. On the right, the LiDAR sensors determine that the polygon is empty, while the radar sensors see and estimate the velocity of a following vehicle to determine lane change feasibility.

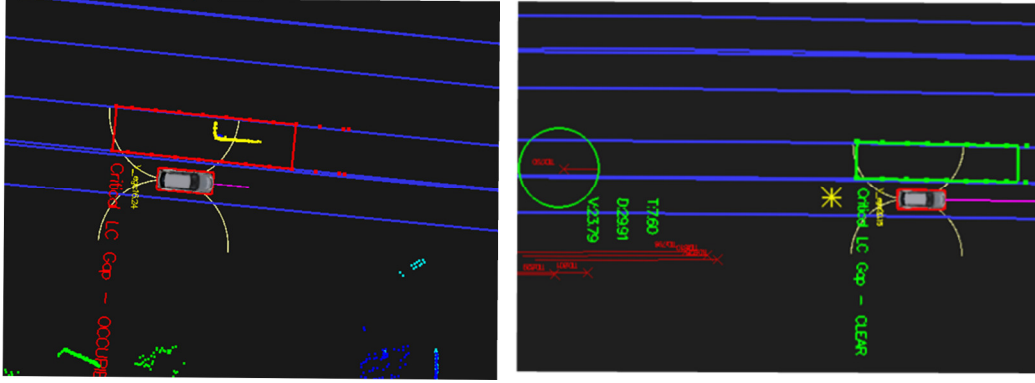


Figure 2. These figures depict a four-lane highway with the two bottom lanes traveling to the right and the two top lanes traveling to the left. The ego vehicle is drawn in the bottom-most lane and is outlined by a red box. *Left*: The red polygon is checked for occupancy, and the LiDAR sensors show a yellow L-shaped outline indicating the presence of a car, which is why the polygon is colored red and lane change is currently infeasible. *Right*: The LiDAR sensors show no car in the polygon, so it is colored green, but there is a radar target in the green circle whose speed is being checked to determine feasibility of lane change.

Merging and Turning. Tests have been performed on the merging/turning situations depicted in Figure 3 and Figure 4. A similar division of labor between LiDAR and radar is used: LiDAR checks polygon(s) near the intersection for occupancy, and radar checks the speed of approaching vehicles for feasibility of merging/turning.

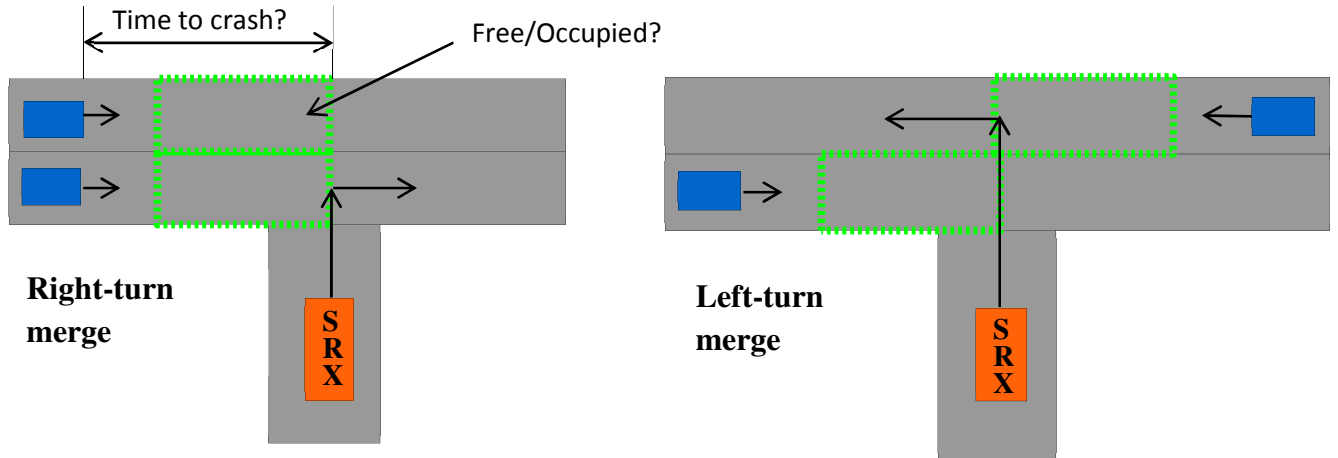


Figure 3. This figure depicts right- and left-turn merging onto a potentially high-speed road from a stop. *Left*: In the right-turn merge case, even though the intent is to turn into the bottom-most lane, we check both lanes in case the SRX swings out a bit into the upper lane. *Right*: In the left-turn merge case, we need both to clear the traffic coming from the left and enter the upper lane comfortably in front of traffic coming from the right.

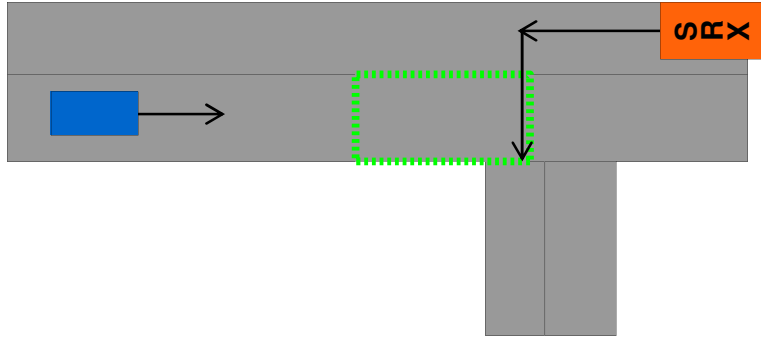


Figure 4. This figure depicts turning left against traffic. The SRX (on the right) wants to turn left at the intersection and has either come to a stop or slowed down in order to do so.

Tests. Both the lane-change and merging/turning algorithms have been tested extensively in real traffic in the Pittsburgh, Pennsylvania area, particularly in and around Cranberry, PA and the highways leading from it to the Pittsburgh International Airport. There has not yet been time to carefully tabulate exact numbers from all logs, but successful lane changes have roughly increased from under 50% to over 85%, and successful right-turn merges from under 25% to about 50% using the algorithms described here, compared to a radar-only or LiDAR-only solution. Failures can result from perception or other factors (e.g., motion planning, control, system latencies), and it is not yet possible definitively to ascribe the failures to their respective causes in every case. However, the general improvement by using case-specific sensor fusion is noteworthy.

Conclusions/Recommendations

This project has resulted in a) a state-of-the-art, relatively computationally efficient pedestrian detection system; and b) a computationally efficient, heuristics-based method of sensor fusion for traffic detection using automotive-grade radar and near-automotive-grade LiDAR. Future work should include:

- Refinement and hardening of the pedestrian detector for deployment and testing in real vehicles and scenarios

- Careful characterization of the traffic detection method results, disambiguating between different failure sources
- Extension of the traffic detection method to other specific cases, such as ramp merging
- Comparison of the case-specific traffic detection results with a comprehensive tracking method to determine an appropriate division of labor between the two methods, which may have complementary strengths/weaknesses

Works Cited

- Cho, H., Rybski, P., & Zhang, W. (2012). Real-time Pedestrian Detection with Deformable Part Models. *IEEE Intelligent Vehicles Symposium*, (pp. 1035-1042).
- Urmson, C., & et al. (2008). Autonomous Driving in Urban Environments: Boss and the Urban Challenge. *Journal of Field Robotics*, 425-466.

Real-time Pedestrian Detection with Deformable Part Models

Hyunggi Cho, Paul E. Rybski, Aharon Bar-Hillel and Wende Zhang

Abstract—We describe a real-time pedestrian detection system intended for use in automotive applications. Our system demonstrates superior detection performance when compared to many state-of-the-art detectors and is able to run at a speed of 14 fps on an Intel Core i7 computer when applied to 640×480 images. Our approach uses an analysis of geometric constraints to efficiently search feature pyramids and increases detection accuracy by using a multiresolution representation of a pedestrian model to detect small pixel-sized pedestrians normally missed by a single representation approach. We have evaluated our system on the Caltech Pedestrian benchmark which is currently the largest publicly available pedestrian dataset at the time of this publication. Our system shows a detection rate of 61% with 1 false positive per image (FPPI) whereas recent other state-of-the-art detectors show a detection rate of 50% ~ 61% under the ‘reasonable’ test scenario (explained later). Furthermore, we also demonstrate the practicality of our system by conducting a series of use case experiments on selected videos of Caltech dataset.

I. INTRODUCTION

Vision-based pedestrian detection is a popular research topic in the computer vision community due to direct application of topics such as visual surveillance [20] and automotive safety [18], [14]. In the past few years, impressive progress has been made, such as found the works of [17], [19], [4], [11], [24], [6], and these insights suggest that they can and should be applied to real world applications. Within the past few years, extremely challenging real-world datasets such as the Caltech Pedestrian [7] and Daimler Pedestrian [8] collections have been introduced. These public datasets allow researchers to evaluate their algorithm’s performance against that of other researchers. We have been making use of these in our research.

Carnegie Mellon University won the 2007 DARPA Urban Challenge with the autonomous vehicle “Boss” [9]. However, in that race, no pedestrians were allowed on the track. While impressive, the technology necessary to win that race is insufficient for operating on real roads. Our recent work has been to focus primarily on developing a real-time pedestrian detection system for automotive safety applications which could be used by both an autonomous vehicle to help it safely navigate roads as well as by a manually-driven vehicle as an early-warning system for the driver. The contributions of this work are threefold.

H. Cho and P. E. Rybski are with ECE and the Robotics Institute, respectively, Carnegie Mellon University, 5000, Forbes Ave., Pittsburgh, PA 15213, USA. {hyunggi, prybski}@cs.cmu.edu

A. Bar-Hillel is with the Advanced Technical Center, General Motors, Hamanofim 11, Herzliya, Israel. aharon.barhillel@gm.com

W. Zhang is with the Electrical and Controls Integration Lab, General Motors R&D, 30500, Mound Rd, Warren, MI 48092, USA. wende.zhang@gm.com

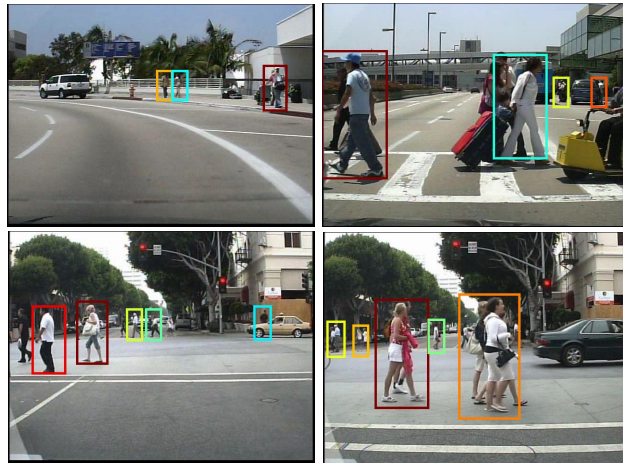


Fig. 1. Typical detection results of our on-board pedestrian detection system. With a multiresolution deformable part-based model, the system can detect pedestrians up to 25m reliably.

The first main contribution is a C implementation of our detection system suitable for real-time operation. Achieving real-time detection speed is by no means a trivial process especially for most of the state-of-the-art detectors which usually use a sliding window method. Our detection system is based on the star-cascade algorithm [10] for a part-based deformable model [11] introduced by Felzenszwalb et al. which is considered as one of the most successful approaches in general object detection. Because of the generality of this method, we can relatively easily apply our system to other relevant categories such as vehicles and bicycles by providing a new training set. In our implementation, we simplify the structure of the star-cascade model for improved speedup and demonstrate a real-time performance with a vehicle-mounted camera.

The second main contribution is a safe geometric constraint analysis based on known camera calibration information for efficient search. Usually, the availability of such information is not assumed in general object detection applications, but for automotive applications where the camera will be mounted in a known position in the vehicle, this information is very useful to exploit. By doing so, we are able to not only accelerate our detection process by searching only relevant image spaces, but we are also able to improve detection accuracy by suppressing a number of potential false positives from irrelevant image space. Based on some assumptions, we analyze the following relationships: ‘pedestrian height in pixels’ vs. ‘distance from a vehicle’ and

‘pedestrians’ foot position in images’ vs. ‘distance from a vehicle’. We propose a simple algorithm for fast search based on these relationship.

Our third main contribution is a quantitative evaluation of our system using public real world datasets. First, we compared our detector’s performance with current state-of-the-art detectors using the Caltech Pedestrian Benchmark. The Caltech dataset is at least two orders of magnitude larger than most exting datasets and provides us a unique and interesting opportunity for in-depth system evaluation. Through a series of well designed experiments, we seek to identify values for key design parameters of our detector such as the optimal number of parts for our deformable part-based model, the optimal number of scales per octave for multiscale detection, etc.

The remainder of this paper is organized as follows. Section II reviews related work on pedestrian detection. Technical implementation details of our detection system are described in Section III and a geometric constraint from known camera calibration information is exploited for an efficient search in Section IV. We describe experimental results using the system in Section V and conclude in Section VI.

II. RELATED WORK

For a comprehensive survey of recent works in vision-based pedestrian detection, refer to [7], [15]. Dollár et al. [7] focuses primarily on the pedestrian detection problem and performs an intensive evaluation of the majority of the state-of-the-art detector algorithms. Gerónimo et al. [15] focuses on pedestrian protection systems for advanced driver assistance systems which utilize tracking, scene geometry, and stereo systems. We review here only important advances for pedestrian detection (not tracking) and describe how these detection approaches can be specially designed for automotive applications.

Historically, one of the first pioneering efforts was the work of Papageorgiou et al. [17] which used Haar wavelet features in combination with a polynomial Support Vector Machine (SVM). They also introduced the first generation pedestrian dataset, known as the ‘MIT Pedestrian Dataset’. Inspired by their work, Viola and Jones [19] brought several important ideas into this field including the use of a new machine learning algorithm (AdaBoost) for automatic feature selection, the use of a cascade structure classifier for efficient detection, and finally the use of an integral image concept for a fast feature computation. Later, the authors demonstrated how to incorporate space-time information into their simple Haar-like wavelet features for moving people detection [20].

The next breakthrough in the detection technology occurred in a feature domain itself for pedestrian detection. Dalal and Trigg [4] showed excellent performance for detecting a human in a static image using dense HOG (Histogram of Oriented Gradient) features with linear SVM. They also introduced the second generation pedestrian dataset, called the ‘INRIA Person Dataset.’ Currently, HOG is considered

to be the most discriminative single feature and is used in nearly all modern detectors in some forms [7].

Detectors that improve upon the performance of HOG utilize a fusion of multiple features and part-based approaches. Wojek and Schiele [23] exploit several combinations of multiple features such as Haar-like features, shapelets, shape context, and HOG features. This approach is extended by Walk et al. [21] by adding local color self-similarity and motion features. Wang et al. [22] combined a texture descriptor based on local binary patterns with HOG. Recently, Dollár et al. [6] provided a simple and uniform framework for integrating multiple feature types including LUV color channels, grayscale, and gradient magnitude quantized by orientation. That group also implemented a near real-time version of this algorithm [5] which makes this method suitable for automotive applications.

Part-based approaches have gained popularity recently mainly because they can handle the various appearances of pedestrians (due to clothing, pose, and occlusion) and can provide a more complex model for pedestrian detection. Mohan et al. [17] use this methodology to divide the human body into four parts: head, legs, left arm, and right arm. Each part detector is trained using a polynomial SVM where outputs are fed into a final classifier after checking geometric plausibility. Mikolajczyk et al. [16] model humans as assemblies of parts that are represented by SIFT-like orientation features. Felzenszwalb et al. [11] demonstrated that their deformable part-based model human detector can outperform many of existing current single-template-based detectors [20], [4], [22]. Based on a variation of HOG features, they introduce a latent SVM formulation for training a part-based model from overall bounding box information without part location labels. Bar-Hillel et al. [1] introduced a new approach for learning part-based human detection through feature synthesis.

Pedestrian detection algorithms have an obvious extension to automotive applications due to the potential for improving safety systems. In this case, the design criterion for a detector might be very different as real-time operation is just important as high detection accuracy. Shashua et al. [18] proposed a part-based representation in a fixed configuration for pedestrian detection. The authors used 13 overlapping parts with HOG-like features and ridge regression to learn a classifier for each part and reported a classification performance of a 93.5% detection rate and a 8% false positive rate. Gavrila and Munder [14] proposed a pipeline using Chamfer matching and several image based verification steps for a stereo camera setup. The classification performance was reported as a 90% detection rate with a 10% false positive rate.

Our group has been making use of the deformable part-based model [11] as part of a tracking-by-detection system that is intended for use with an autonomous vehicle as well as an early-warning system for driver safety. We have demonstrated the algorithmic viability of these methods in off-line analysis [2], [3] and in this work we seek to demonstrate how these approaches can be implemented in real-time.

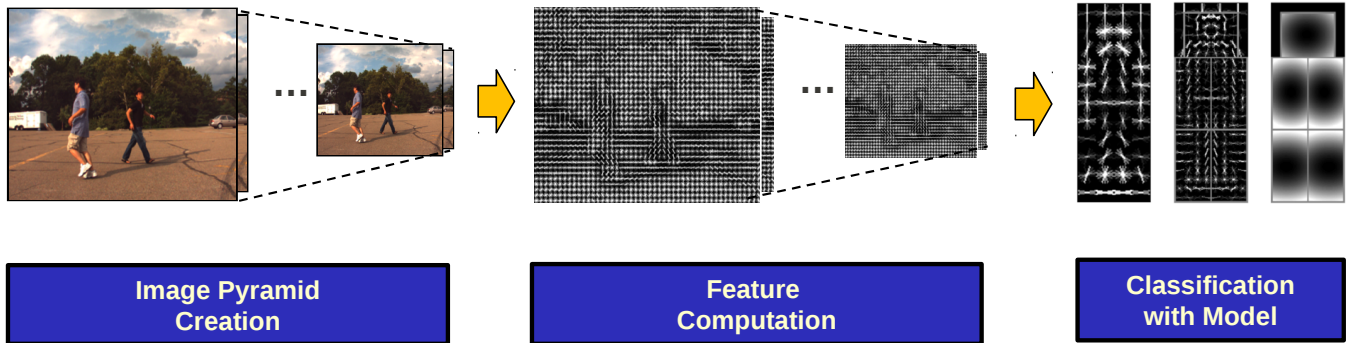


Fig. 2. Illustration of a procedure for a multiscale pedestrian detection. The main computational bottlenecks are feature computation for densely sampled image pyramid (HOG features in our case) and sliding-window classification (a large number of cross-correlation operations with several parts in our case).

III. IMPLEMENTATION OF DEFORMABLE PART-BASED MODELS

Our approach is based on the deformable part-based HOG model [11]. We have engineered this algorithm in an attempt to optimize it to be the core of a real-time pedestrian detection for automotive applications. To achieve this, we re-implemented the algorithm (originally based in MATLAB/C) using C and optimized several subsystems in an attempt to improve the algorithms overall speed. The remainder of this section discusses the main reasons for choosing the deformable part-based model and describes the important details of our implementation that were key to its fast performance.

A. Why Deformable Part-Based Model ?

Real-time onboard pedestrian detection from a moving vehicle is a challenging task especially when using a monocular camera as the primary (and sole) sensor. Although a number of approaches have been proposed, there are still a number of aspects of this problem that are difficult and can be considered to be as of yet “unsolved.” That being said, the accuracy and computational performance of these detectors are improving with each successive publication. We performed a broad survey across the current state-of-the-art pedestrian detectors and ended up with the following methods [22], [11], [5], [1] as candidates for our implementation.

We opted to implement a detector based on the work of Felzenszwalb et al. [11] for the following reasons. First, this method provides an elegant mechanism for handling a wide range of intra-class variability (i.e., various poses of pedestrians) by having multiple submodels (so-called “components”). Furthermore, the deformable part-based model exploits dynamic configurations of its parts which allows for a wider variation in object appearances to be accepted. The importance of this aspect is well illustrated in a recent survey [7]. Secondly, this method has a well-designed learning technique called “latent SVM” which not only can handle a large training set effectively, but also can learn a part-based model without requiring information about exact part labels. Finally, this approach utilizes an efficient detection

TABLE I
PROFILING RESULTS FOR ALL IMPLEMENTATIONS (UNIT:MS)

Algorithm Name	Computer I Intel Core2 Duo P8800@2.66GHz 2GB RAM		Computer II Intel Core i7 2920XM@2.5GHz 16GB RAM	
	Matlab star-cascade	C star-cascade	Matlab star-cascade	C star-cascade
HOG Feature Computation	1145	300	840	80
Sliding Window Classifier	560	320	300	100
Non-Maximal Suppression	24	10	24	5

mechanism called the star-cascade algorithm [10] which makes it suitable for real-time applications.

B. Implementation Details

Our implementation follows a standard procedure for processing the image data that consists of first creating a densely sampled image pyramid, computing features at each scale, performing classification at all possible locations, and finally performing non-maximal suppression to generate the final set of bounding boxes. The key steps for this process are illustrated in Figure 2. The key factors that we found for producing the best results for our algorithm include densely sampling the image pyramid, computationally intensive (somewhat tedious) classification in the search space, and the use of discriminative features such as HOG.

Our final goal in terms of speed performance was a speed of 10fps on 640×480 resolution images. To understand where our efforts would need to be applied in the implementation of this algorithm, we profiled the performance of the original MATLAB/C based detectors. These detectors included the `voc-release3` and `voc-release4` with a star-cascade algorithm¹ [13]. The profiling was performed using two evaluation computers that included an Intel Core2

¹<http://www.cs.brown.edu/~pff/latent/>, accessed on May 2011

Duo P8800@2.66GHz with 2GB RAM, labeled computer I, and an Intel Core i7 2920XM@2.5GHz with 16GB RAM, labeled computer II. For 640×480 images and 10 scales per octave, `voc-release3` (1 component with 6 parts, multi-threaded version) demonstrated a performance of 0.5 fps and `voc-release4` algorithm with a star-cascade algorithm (1 component with 8 parts, single-threaded version) demonstrated a performance of 0.6 fps on computer I. The details of the profiling result is shown in Table I. As can be seen in this Table, the two main bottlenecks are the computation of HOG features for densely sampled image pyramids and the sliding-window classification. Since the MATLAB/C based detector already uses compiled C code (MEX functions in MATLAB) we needed to approach these bottlenecks by developing parallel (multi-threaded) implementation of the HOG feature computation functions as well as the classification functions. The key details to this implementation are described below:

Feature computation: Given an input image, the first step of the pipeline is to build a densely sampled image pyramid and compute the corresponding feature pyramid. For an image pyramid, we employed the image resize function of the OpenCV 2.0 library. For the HOG feature pyramid, which was the first computational bottleneck, we refactored the algorithm to make use of the `pthread` library. This was possible because the computational process to generate each level of the pyramid is independent of all the others. This solution allowed us to speed the algorithm up by 1 order of magnitude.

Sliding-window classification: For the `voc-release3` algorithm which is based on star-structured models (1 root filter and 6 part filters), this process at a certain pyramid level corresponds to 7 times cross-correlations between each model and that level of feature pyramid and then 6 times generalized distance transforms [12] to combine all part filter responses. In practice, a large number of cross-correlation is the main bottleneck for this step. For this, the original method provides a parallelized correlation scheme with a numerical linear algebra enhanced function. We ported their MEX functions into our implementation. Whereas, `voc-release4` with a star-cascade algorithm provides an efficient way for classification by evaluating parts in a cascaded fashion with a sequence of thresholds. For implementation of this idea, `voc-release4` uses $2(n+1)$ models for a model with $n+1$ parts, where the first $n+1$ models are obtained by sequentially adding parts with simplified appearance models (PCA version of original HOG feature) for faster computation and second $n+1$ models are obtained by sequentially adding parts with its full feature. Our implementation is a little bit different in that we just use $n+1$ cascade models with full HOG feature for ease of implementation and we parallelize the process using `pthread` library. By doing this, we achieve 2X-4X speed improvement.

Non-maximal suppression: After sliding-window classification, we usually get multiple overlapping bounding boxes from detections nearby locations and scales. Non-maximal suppression (NMS) is used for eliminating those repeated detections. Currently, two dominant NMS approaches have

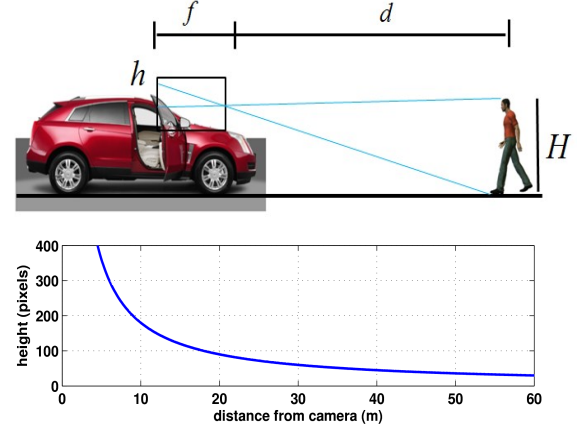


Fig. 3. Geometric constraints analysis. (a) Scene geometry. (b) Pixel height h as a function of distance d .

been widely used: mean shift mode estimation introduced in [4] and pairwise max suppression introduced in [11]. We implemented the pairwise max suppression scheme due to its simplicity and efficiency.

IV. GEOMETRY ANALYSIS

The use of scene geometry enables two primary benefits for pedestrian detection: efficient and accurate detection. Efficiency can be achieved by searching only geometrically valid regions in the image space and accuracy can be achieved by suppressing a number of potential false positives from the irrelevant image space. Many automotive applications therefore exploit geometric constraints from the known camera calibration information with some assumptions [18], [14]. Here, we propose a simple and efficient search algorithm for our pedestrian detector by analyzing the following relationships: ‘pedestrian height in pixels’ vs. ‘distance from a vehicle’ and ‘pedestrians’ foot position in images’ vs. ‘distance from a vehicle’. For analysis for these relationships, we assume that 1) ground plane is flat, 2) pedestrians rest on the ground plane, and 3) localization of bounding boxes is accurate (especially for bottom lines of bounding boxes).

Pedestrian height in images: Since Caltech dataset was collected with a typical automotive camera settings (640×480 resolution, 27° vertical field of view, and fixed focal length (f) at 7.5mm), we can compute a pixel height h of a pedestrian in a image using a perspective projection model. From Figure 3(a), we can easily draw the relationship between a pixel height (h) and distance (d) as $h \approx H f_p / d$, where H is the true pedestrian height and f_p is the focal length expressed in pixels. This relationship for Caltech dataset is shown in Figure 3(b) assuming $H = 1.8m$.

Depth computation: In general, estimating depth information from monocular camera images is impossible. However, if all assumptions we mentioned above are true, we can compute a range from the camera to a pedestrian. Here,

we avoid the derivation of this relationship due to space limitation. The derivation is provided in our project website².

We can search an input image efficiently by computing depth information first according to a current y -coordinate of searching step and then selecting geometrically plausible scales of a feature pyramid so that a fixed size pedestrian model can detect real pedestrians. By this search scheme, we can achieve 2X-3X speed improvement.

V. EXPERIMENTAL RESULTS

To quantitatively evaluate our pedestrian detection system, we analyzed its performance on various real-world datasets. To evaluate detection performance, we used the Caltech Pedestrian Dataset [7] which offers an opportunity to exploit many different aspects of model training thanks to its large number of positive samples. The (annotated) dataset corresponds to approximately 2.5 hours of video captured at 30fps. The dataset is segmented into 11 sessions, 6 of which were used for training (S0~S5) and the rest sessions were used for testing (S6~S10). We performed a set of experiments to identify the key design parameters for the deformable part-based model. Using the model trained with the optimal parameters from the first set of experiments, we compare our performance with other state-of-the-art detectors using the Caltech Benchmark. This evaluation framework required us to upscale the input images by a factor of 2 (with a final resolution of 1280×960) for the first and second experiments. For subsequent experiments, we fabricated our own set of experiments (what we call 'automotive') to better represent data that would be seen from real-time automotive applications.

A. System Design Parameters

Number of Training Samples: The Caltech dataset contains approximately 2,300 unique pedestrians out of approximately 347,000 total instances of pedestrians. Because of this high level of redundancy and correlation in the sample images, we had to first determine the best sampling scheme to obtain a statistically valid set of training images. We trained 7 models with a standard set of parameters (1 component 6 parts) where each model trained with a training set (S0~S5) consisting of images selected from the full data set at a different sampling frequency. We trained models for each of the following frequencies: 1, 10, 20, 30, 40, 50, and 60. We ran each model on the test set (S6~S10) and evaluated them using the same evaluation code provided from [7] (ver. 3.0.0). The results of these experiments are shown in Table II, where we use log-average miss rate (LAMR) to represent detector performance by a single reference value. The LAMR is computed by averaging miss rate at nine FPPI rates evenly spaced in log-space in the range 10^{-2} to 10^0 according to [7].

Although there is no large difference in performance between the settings, using every 30th training images (i.e., one training image per one second) in this case gives us the best performance. We decided to use this setting throughout

TABLE II
DIFFERENT SAMPLING SCHEMES FOR MODEL TRAINING

Sampling Scheme	No. of Pos. Samples	LAMR (%)
All frames	65570	56
Every 10th	6557	56
Every 20th	3261	55
Every 30th	2198	54
Every 40th	1629	55
Every 50th	1280	56
Every 60th	1088	58

the remainder of our experiments. This result has implication in the generation of ground truth data as annotating every 30th frame of a dataset is far less expensive than having to annotate every frame.

Number of Parts: The standard number of parts for the deformable part-based model is 6 and 8 for `voc-release3` and `voc-release4`, respectively. However, the optimal number of parts depends on the variability of an object class and may be significantly different between classes. Our second experiment was designed to identify the optimal number of parts required for the pedestrian models that would be used for the automotive application. Once again, we trained 7 models with different number of parts (2~8) using the exact same training set and tested them on the testing set. As shown in Table III, using the standard part number of 6 shows the best performance.

TABLE III
DIFFERENT NUMBER OF PARTS

No. of Parts	LAMR (%)
2	58
3	55
4	55
5	56
6	54
7	56
8	56

In general, when selecting the number of parts, say for `voc-release3`, it might be better to use 3 or 4 parts for a faster detection rate. However, for the star-cascade algorithm using 6 to 8 parts seems reasonable due to the cascaded structure of its classifier. To balance efficiency and accuracy, we decided to use 6 parts for our evaluations.

Number of Scales Per Octave: Typically, modern pedestrian detectors use two or three octaves and sample 8-14 scales per octave for accurate detection. Since the computational requirements for feature computation in this image pyramid can be significant, new methods of handling this issue must be identified. Recently, Dollár et al. [5] proposed a technique to avoid constructing such a feature pyramid by approximating feature responses at certain scales by computing feature responses at a single scale. They demonstrated that the approximation is accurate within an entire scale octave. Here, we are interested in primarily looking at performance

²www.ece.cmu.edu/~hyungic/pedestrian.html

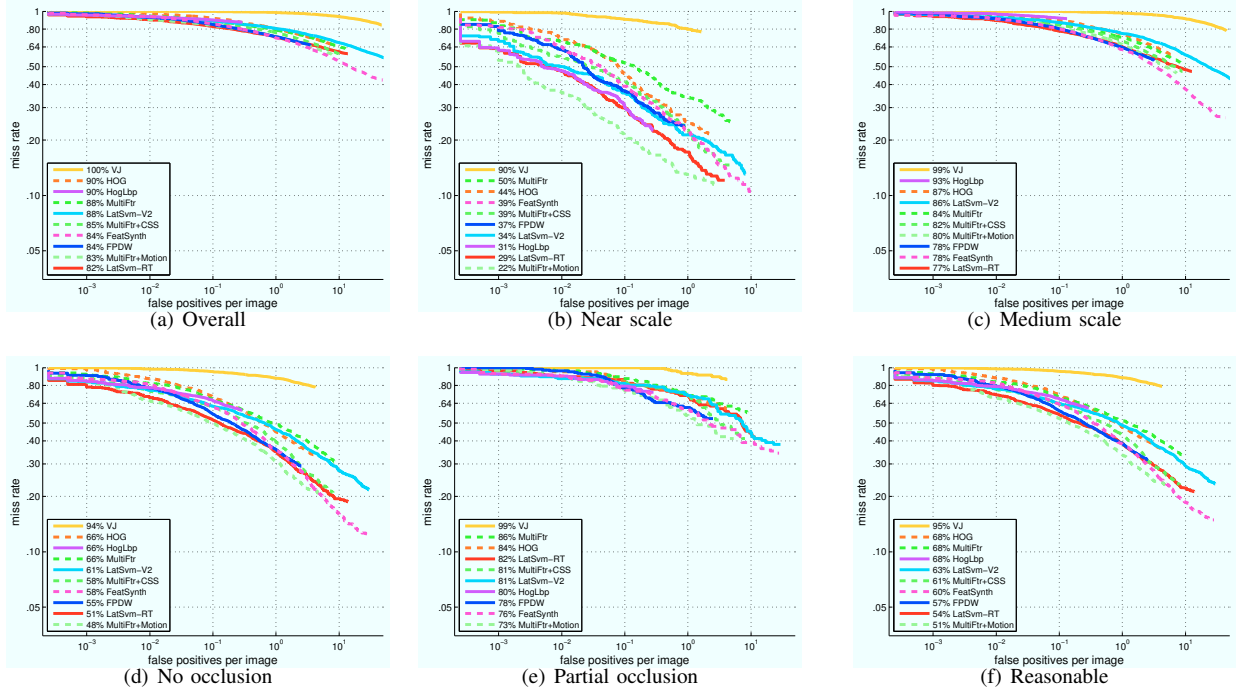


Fig. 4. Evaluation results using the same criteria in [7]. (a) Overall performance on all annotated pedestrian images. (b) Performance on unoccluded pedestrians over 80 pixels in height. (c) Performance on unoccluded pedestrians between 30-80 pixels in height. (d) Performance on unoccluded pedestrians over 50 pixels in height. (e) Same as (d) but with partial occlusion. (f) Performance on pedestrians at least 50 pixels in height under no or partial occlusion.

differences depending on different number of scales and looking for a specific optimal solution for our configuration. We tested 7 different settings and the results are shown in Table IV. We found that even 4 scales per octave shows a marked improvement over 2 scales per octave in accuracy. We decided to use 4 scales per octave as a trade-off between accuracy and performance.

TABLE IV
DIFFERENT NUMBER OF SCALES PER OCTAVE

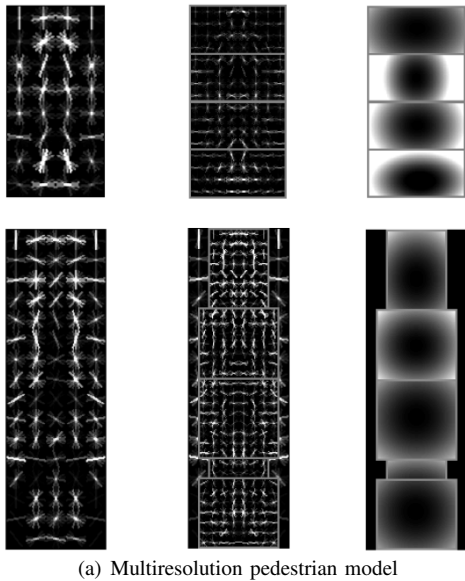
Scales / Octave	No. of Levels	LAMR (%)
2	10	61
4	19	56
6	28	56
8	37	55
10	46	54
12	56	54
14	65	53

B. Evaluation with Caltech Benchmark

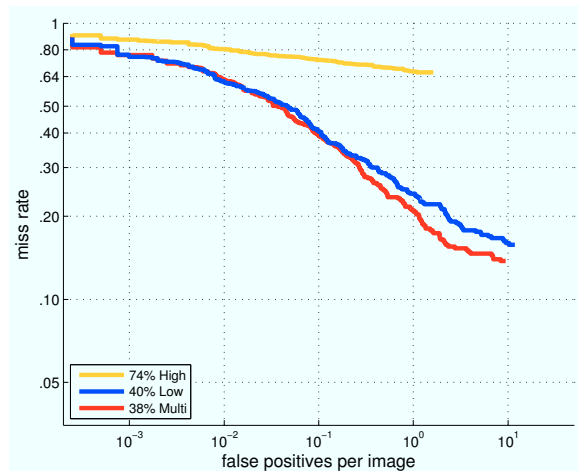
Using a model parameterized with the values identified in our previous experiments, we compared the performance of our detector with those of other state-of-the-art detectors using the Caltech evaluation framework. This framework provides a flexible way to evaluate detectors under various conditions on multiple datasets. We evaluated the performance of the algorithms under six conditions on the testing data in the Caltech dataset. Although the framework provides detection results from sixteen pre-trained detectors to help

provide consistent comparison, we intentionally selected 10 of the most relevant top-performing detectors to increase the clarity of the resulting graphs. The criterion we used to select the algorithms was detection accuracy and runtime. One exception is the “MultiFtr+Motion” algorithm, which shows very slow detection speed (estimated around 30 seconds for a 640×480 image on our evaluation computer II). But we included the method since it is the only method which uses motion features in our list and because it shows the best performance in most cases. We named our algorithm “LatSvm-RT”. Following the naming convention for the Caltech benchmark we described the other detectors as follows: VJ [19], HOG [4], HogLbp [22], MultiFtr [23], LatSvm-V2 [11], MultiFtr+CSS [21], FeatSynth [1], FPDW [5], MultiFtr+Motion [21]. All detectors except MultiFtr+CSS [21] and MultiFtr+Motion [21] (which use their own dataset called ‘TUD-MP’) were trained with the ‘INRIA Person Dataset [4].’ Note that the LatSvm-V2 is our baseline detector, *voc-release3*. Our results are shown in Figure 4. Following the example found in [7], we plot miss rate vs. false positives per image (FPPI) and use the log-average miss rate as a single comparison value for the overall performance. The entries are ordered in the legend from the worst log-average miss rate to the best.

In general, our detector shows very good performance except for the “partial occlusion” test case. Figure 4(a) illustrates algorithmic performance on the entire test set which contains pedestrians from 20 pixels in height which are in general fairly poor. Results for near and medium scale unoccluded pedestrians, corresponding to heights of at least



(a) Multiresolution pedestrian model



(b) Automotive

Fig. 5. Evaluation results using the *automotive* criterion. (a) Multiresolution pedestrian model: high-resolution (128×40) with 475 positive samples and row-resolution (64×32) with 2430 positive samples. (b) Performance on unoccluded pedestrians over 70 pixels tall.

80 pixels and 30-80 pixels, respectively, are shown in Figure 4(b) and Figure 4(c). In each case, our detector shows the second best or best performance with a log-average miss rate of 29% and 77%, respectively. As can be seen in Figure 4(d) and Figure 4(e), performance drops significantly under partial occlusion, leading to a dramatic decrease of a log-average miss rate from 51% to 82% for our case. This is quite disappointing since our classification algorithm can not handle partial occlusions even though the part-based model itself has a rich representation for partial occlusion. We will seek to address this through the development of a new classification algorithm for partial occlusion handling in our future work. Finally, Figure 4(f) shows performance on pedestrians over 50 pixels tall under no or partial occlusion. Here our detector shows second best performance with a log-average miss rate of 54%.

C. Real-Time Evaluation

Because the goal of our research is to develop pedestrian detection algorithms that are suitable for the rigors of deployment on an autonomous vehicle, we developed a new test scenario that we call ‘*automotive*’ to show the performance of our detector in this context. We define real-time operation to be 10fps@ 640×480 . For the automotive experiment, the ‘*reasonable*’ scenario in the Caltech benchmark does not fit to our needs for several reasons. First, it uses detection results from upscaled input images to compare performances. However, it is almost impossible to process these images at 10Hz without an unreasonable amount of computational power. Second, we need to adjust ground truth information according to our working distance. In the ‘*automotive*’ setting, we can only include unoccluded pedestrians within 25m from a vehicle (corresponds to 70 pixels in height) in the ground truth. To increase the working range of our

detector while maintaining a high detection rate, we trained a multiresolution pedestrian model (shown in Figure 5(a)). We intentionally set the height of our low-resolution model as 64 pixels so that our system can detect pedestrians up to 25m reliably. We run our detector on the Caltech testset without upscaling (i.e., 640×480 images) with the efficient search scheme discussed in Section IV. We achieve a real-time operation at 14fps with a log-average miss rate of 38%. We believe our work is a meaningful progress of state-of-the-art pedestrian detectors. The performance is shown in Figure 5(b) and some qualitative results achieved on the Caltech testset are shown in Figure 6.

VI. CONCLUSIONS AND FUTURE WORK

This paper presents a real-time pedestrian detection system for intelligent vehicles using a deformable part-based model. For support real-time operation, we implement two different versions of Felzenszwalb et al.’s object detection systems (a baseline [11] and a star-cascade method [10]). Scene geometry from known camera calibration information is utilized to search a feature pyramid more efficiently. For better detection accuracy, a multiresolution pedestrian model is used for detecting small (pixel-sized) pedestrians as well as normally-sized ones. Using the Caltech Pedestrian Dataset, we quantitatively evaluated our detection system with a series of experiments and showed real-time operation (14fps@ 640×480 images) while maintaining the state-of-the-art detection accuracy (80% detection rate with 1 FPPI) under our test scenario called ‘*automotive*’. As part of our future work we want to develop a partial occlusion handling algorithm to increase detection accuracy.

VII. ACKNOWLEDGMENTS

This project was funded by General Motors through the General Motors-Carnegie Mellon Autonomous Driving Col-



Fig. 6. Qualitative detection results on the Caltech testset. The first and second row shows correct pedestrian detections in various scenarios. The third row shows typical false positives.

laborative Research Lab.

REFERENCES

- [1] A. Bar-Hillel, D. Levi, E. Krupka, and C. Goldberg. Part-based feature synthesis for human detection. In *ECCV*, 2010.
- [2] H. Cho, P. Rybski, and W. Zhang. Vision-based bicycle detection and tracking using a deformable part model and an ekf algorithm. In *ITSC*, 2010.
- [3] H. Cho, P. Rybski, and W. Zhang. Vision-based 3d bicycle tracking using deformable part model and interacting multiple model filter. In *ICRA*, 2011.
- [4] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005.
- [5] P. Dollár, S. Belongie, and P. Perona. The fastest pedestrian detector in the west. In *BMVC*, 2010.
- [6] P. Dollár, Z. Tu, P. Perona, and S. Belongie. Integral channel features. In *BMVC*, 2009.
- [7] P. Dollár, C. Wojek, B. Schiele, and P. Perona. Pedestrian detection: An evaluation of the state of the art. *PAMI*, 99(PrePrints), 2011.
- [8] M. Enzweiler and D. Gavrilu. Monocular pedestrian detection: Survey and experiments. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 31:2179–2195, 2008.
- [9] C. U. et al. Autonomous driving in urban environments: Boss and the urban challenge. *Journal of Field Robotics, Special Issue on the 2007 DARPA Urban Challenge, Part I*, 25(8):425–466, 2008.
- [10] P. Felzenszwalb, R. Girshick, and D. McAllester. Cascade object detection with deformable part models. In *CVPR*, 2010.
- [11] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 99, 2010.
- [12] P. Felzenszwalb and D. Huttenlocher. Distance transforms of sampled functions. *Cornell Computing and Information Science Technical Report TR2004-1963*, 2004.
- [13] P. F. Felzenszwalb, R. B. Girshick, and D. McAllester. Discriminatively trained deformable part models, release 4. <http://people.cs.uchicago.edu/~pff/latent-release4/>.
- [14] D. M. Gavrila and S. Munder. Multi-cue pedestrian detection and tracking from a moving vehicle. *International Journal of Computer Vision*, 73:41–59, 2007.
- [15] D. Geronimo, A. Lopez, and T. G. A. Sappa. Survey of pedestrian detection for advanced driver assistance systems. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 32, 2010.
- [16] K. Mikolajczyk, C. Schmid, and A. Zisserman. Human detection based on a probabilistic assembly of robust part detectors. In *ECCV*, 2004.
- [17] C. Papageorgiou and T. Poggio. A trainable system for object detection. *International Journal of Computer Vision*, 38(1):15–33, 2000.
- [18] A. Shashua, Y. Gdalyahu, and G. Hayun. Pedestrian detection for driving assistance systems: Single-frame classification and system level performance. In *IEEE IV*, pages 13–18, 2004.
- [19] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *CVPR*, 2001.
- [20] P. Viola, M. J. Jones, and D. Snow. Detecting pedestrians using patterns of motion and appearance. *International Journal of Computer Vision*, 63(2):153–161, 2005.
- [21] S. Walk, N. Majer, K. Schindler, and B. Schiele. New features and insights for pedestrian detection. In *CVPR*, 2010.
- [22] X. Wang, T. X. Han, and S. Yan. An hog-lbp human detector with partial occlusion handling. In *ICCV*, 2009.
- [23] C. Wojek and B. Schiele. A performance evaluation of single and multi-feature people detection. In *DAGM Symposium Pattern Recognition*, 2008.
- [24] C. Wojek, S. Walk, and B. Schiele. Multi-cue onboard pedestrian detection. In *CVPR*, 2009.