

Unsupervised Processing of Vehicle Appearance for Automatic Understanding in Traffic Surveillance

Jakub Sochor¹, Adam Herout

Graph@FIT, Brno University of Technology, Czech Republic
Email: isochor,herout@fit.vutbr.cz

Abstract—This paper deals with unsupervised collection of information from traffic surveillance video streams. Deployment of usable traffic surveillance systems requires minimizing of efforts per installed camera – our goal is to enroll a new view on the street without any human operator input. We propose a method of automatically collecting vehicle samples from surveillance cameras, analyze their appearance and fully automatically collect a fine-grained dataset. This dataset can be used in multiple ways, we are explicitly showcasing the following ones: fine-grained recognition of vehicles and camera calibration including the scale. The experiments show that based on the automatically collected data, make&model vehicle recognition in the wild can be done accurately: average precision 0.890. The camera scale calibration (directly enabling automatic speed and size measurement) is twice as precise as the previous existing method. Our work leads to automatic collection of traffic statistics without the costly need for manual calibration or make&model annotation of vehicle samples. Unlike most previous approaches, our method is not limited to a small range of viewpoints (such as eye-level cameras shots), which is crucial for surveillance applications.

I. INTRODUCTION

We are dealing with automatic and accurate visual traffic surveillance. In this particular work, we aim at accurate and automatic calibration of camera including scene scale, and at fine-grained recognition of passing vehicles. We show that both these problems can be successfully approached by analyzing observed traffic in an unsupervised manner. Instead of providing the system with a large degree of supervision – performing camera calibration procedures and manually annotating sufficiently large datasets of vehicle samples – we let it collect as much information as possible by itself and provide small amount of annotation ex post, as the final bit needed for its understanding. Figure 1 shows examples of three groups of vehicles obtained fully automatically on different cameras. We show in this paper that once having these exact make&model examples, it is possible to calibrate the camera very accurately including scale, and to train specific car detectors. A crucial point is to distinguish between very subtle visual differences, and at the same time, we want the system to include in one class samples of an identical car model under varying lighting conditions, in different environments, different in their color, etc.

Recognition of exact make and model of cars has attracted

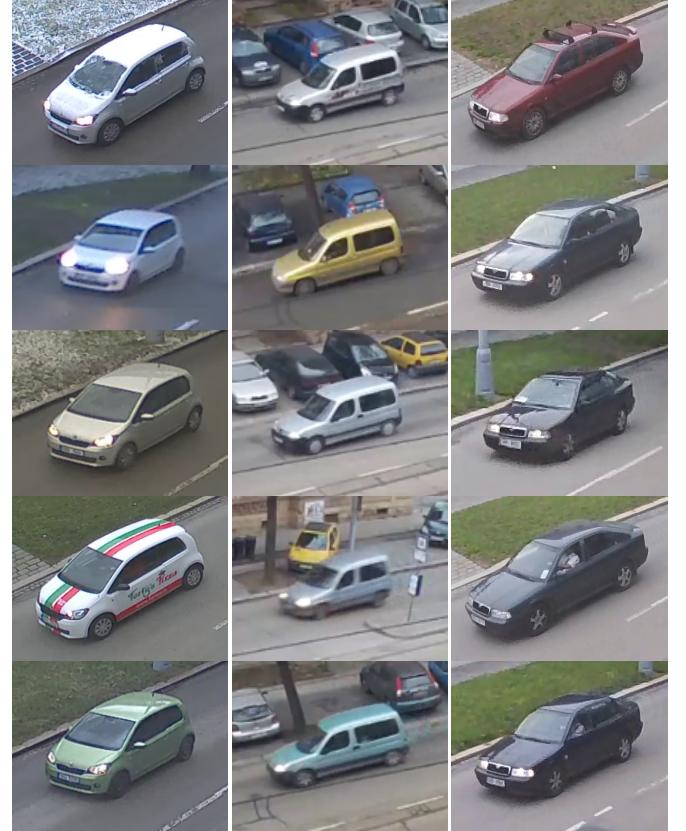


Fig. 1. Examples of groups of vehicles obtained from different cameras at different locations. The groups were obtained from acquired dataset of 472 k vehicles (3 images per vehicle) from 7 different camera positions at multiple locations in multiple countries.

more attention recently. Most approaches are limited to a single viewpoint or a narrow range of viewpoints [1], [2], [3], [4], [5]. Only recently, Hsiao et al. [6] proposed a make and model recognizer operating from freely displaced viewpoints. They train on rendered data (still limited to eye-level views) of eight different Honda models. Similar achievements are still limited by training and evaluation data – the sizes of available datasets are rather modest. Collecting datasets of real-world samples for fine-grained car classification is a difficult task. Stark et al. [7] collected a dataset of 14 categories (1904 samples in total) and they comment: “We note that the data set is heavily biased w.r.t. viewpoints, which reflects the availability of images we encountered during data collection. It proved almost impossible to collect more than a handful of images

¹Jakub Sochor is a Brno Ph.D. Talent Scholarship Holder — Funded by the Brno City Municipality. This work was supported by TACR project “RODOS”, TE01020155 and by the research CEZMSMT project “IT4F”, CZ 1.05/1.1.00/02.0070.

for certain combinations of car-type and viewpoint.” Krause et al. [8] mention other difficulties they faced when collecting a larger dataset (16,185 images, 197 classes). They crawl popular car websites and analyze the car models. Based on the list of models they obtained images by using Flickr, Google, and Bing, and annotated/verified the set manually. Among the difficulties, they mention: “*The class list of cars changes on a yearly basis, the appearance of some models of cars remains constant from year to year, and typical car websites may even list cars which differ only in terms of non-visual features, e.g. their engine, as separate classes.*” The viewpoints are also distributed very unevenly – labeled images of cars posted online by fans are typically taken from a limited range of viewpoints. Lin et al. [9] collected their own dataset of 300 images of 30 car models, because they missed landmark annotation in the previously mentioned sets.

When it comes to fine-grained vehicle classification, many approaches are limited to frontal or rear viewpoint and they are based on detection of the license plate for ROI extraction [1], [2], [3], [4]. Authors of these papers are using different schemes for extracting the feature vectors and for the classification itself. Stark et al. [7] use fine-grained categorization of cars in order to obtain metric information and obtain a rough estimate of depth information for single images (containing cars in usable poses). Krause et al. [10] reconstruct 3D models of cars by wide baseline stereo and combine the recognized 3D model with 2D representations (locality constrained linear coding LLC [11] and BubbleBank [12]) for classifying the vehicles. Another approach proposed by Prokaj and Medioni [13] is based on pose estimation and it is able to handle any viewpoint. The authors suggest to use 3D models of vehicles, fit them to the recognized pose, project them to 2D and use SIFT-like features for the comparison of the vehicles.

Existing systems for camera calibration including scale often involve physical measurements in the scene of interest [14], known length of lane marking [15], or marking dimensions [16], known camera position [15], [17] or average vehicle size [18]. A novel state-of-the art approach proposed by Dubská et al. [19] was published recently. This approach uses three-dimensional bounding boxes and mean dimensions of vehicles; therefore, all vehicles are treated the same way and assumed they will have identical dimensions. In contrast to this approach, we propose camera scale calibration which treats groups of vehicles separately and it assumes that only vehicles of the same type have the same dimensions.

We propose to collect fine-grained vehicle classes in a totally unsupervised manner. We process continuous video from a static surveillance camera and track the passing cars. From each tracked car we obtain several (3-4) representative images at predefined locations. These groups of images are analyzed pairwise by a dissimilarity cost function and clustered into clusters of common appearance. The clustering can be done in a conservative way, so that only samples with very low dissimilarity fall into one cluster. These clusters are annotated manually with only little human effort. Along with the make and model annotation, also the viewpoint direction and scale is automatically obtained for each of the samples. For static cameras, the viewpoint direction and scale can be derived from vanishing points that can be fully automatically extracted for a roadside surveillance camera [19].

Clusters obtained by our approach can be matched from camera to camera or between different time spans of observation. Matching clusters from a single camera at different times is important for processing over time. The ability to construct clusters from limited time spans and to join them together later allows for continuous processing of input cameras and for gradually growing the collected dataset. The between-camera matching is used for transferring manual annotations made for one camera to other views so that each make&model cluster contains samples from various viewpoints and of many individual instances of the given car at different locations. The between-camera matching can also transfer camera scales and our algorithm can thus be used for precise automatic camera calibration, including the scale. The accuracy of the clustering and camera scale inference is measured experimentally and reported in Section III.

The main contributions of this paper are the following:

- We propose a new unsupervised approach to traffic understanding and camera calibration. It is specific for traffic surveillance – dataset creation, camera calibration, size+speed measurement, etc.
- Our approach beats the state-of-the-art in traffic camera calibration including scale [19], leading to twice as accurate length and speed measurement.
- Our approach leads to unsupervised traffic dataset collection from all viewpoint angles imaginable in traffic surveillance. This contrasts to majority of available traffic datasets, focusing on eye-level viewpoints with almost constant Yaw angle [20], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31].
- We are making the collected datasets (472k of vehicles, 3 images each, various viewpoints relevant for traffic surveillance) and the source codes publicly available for future development and benchmarking¹.

The structure of the text is following: First, we describe into a detail the algorithms for automatic grouping and matching of vehicles in Section II. The used method for scene scale inference is also described in this section. Then, in Section III, we evaluate our algorithms in numerous scenarios. Finally, we conclude our work in Section IV.

II. UNSUPERVISED SIMILARITY-BASED VEHICLE GROUPING

The images of vehicles for grouping are obtained automatically from traffic surveillance cameras. 3D bounding boxes are constructed around the cars [19], Figure 3. Samples of vehicles are acquired when the center of their base is close to manually predefined positions (see Figure 2). The set of such designated positions will be denoted as \mathbb{P} and an individual position as $\alpha \in \mathbb{P}$. The dimensions and the mask of the vehicle foreground [32], [33] are also acquired for each of the samples. Only vehicles manifested in all of the designated positions α are further processed. In the following text, v denotes a vehicle from the set of vehicles V ; the number of vehicles is denoted $N = |V|$.

Our goal of the following effort is to automatically obtain groups of vehicles with the same make&model (potentially

¹The dataset and source codes are available at medusa.fit.vutbr.cz/traffic/



Fig. 2. Acquisition of vehicle samples. The surveillance camera tracks moving vehicles (green line represents one vehicle track). Several (typically 3) image locations are manually designated (red crosses) and vehicle locations corresponding to them are selected based on distance of the center of the 3D vehicle base location.

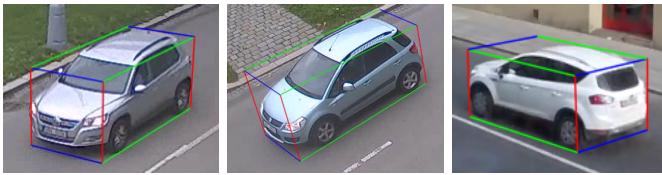


Fig. 3. Examples of three-dimensional bounding boxes constructed according to [19]. Lengths, widths, and heights are denoted by green, blue, and red color.



Fig. 4. A typical example of surveillance camera where vehicle identification or recognition based on SIFT matching fails, as totally different SIFTS are detected at different positions.

also submodel and model year). The computer vision approach of first choice is to extract keypoints and their descriptors and match the images based on such visual keywords. However, Figure 4 shows that in the given image quality and resolution, this approach would not lead to desired outcomes – the SIFTS (other available keypoint detectors behave similarly) are not stable and do not respond to the identical object stably. To deal with this problem, we propose to use edge-based similarity for vehicles as the edges are present on non-textured vehicles even in low resolutions. See Figure 11 for the resulting extracted vehicles.

For each vehicle v and position α , edge points $F_v^\alpha = \{(x_1, y_1), \dots, (x_n, y_n)\}$ are detected by an edge detector [34], [35] and their orientation $O_v^\alpha : F_v^\alpha \mapsto \mathbb{R}$ is also computed. We define a dissimilarity cost between two vehicles $u, v \in V$ in the following way. First, edge points F_v^α are aligned to edge points F_u^α by ICP [36] in order to reduce small movements caused by taking the image at a slightly different position. The transformation aligning points F_v^α to F_u^α obtained by ICP is denoted as ϕ . We propose to use dissimilarity cost d_{uv} , Eq.

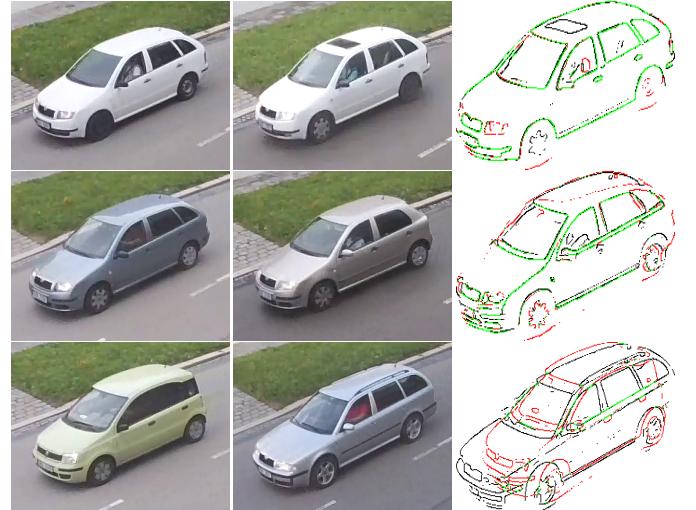


Fig. 5. Dissimilarity cost between vehicle samples. **left:** reference vehicle, **middle:** test vehicle, **right:** edge dissimilarity visualization (black: test vehicle, green: pixels of reference vehicle with a match, red: pixels of reference vehicle without a match).

(1), which is inspired by Chamfer Matching [37]. Threshold λ is (in our experiments) set to 1 and $\triangleleft[\cdot, \cdot]$ denotes the unsigned difference between two unoriented angles from range $(0, \pi)$; therefore, the maximal angular difference is $\frac{\pi}{2}$.

$$d_{uv} = \sum_{\alpha \in \mathbb{P}} d_{uv}^\alpha \quad (1)$$

$$d_{uv}^\alpha = \frac{\sum_{p_u^\alpha \in F_u^\alpha} d(p_u^\alpha, \arg \min_{p_v^\alpha \in \phi(F_v^\alpha)} \|p_u^\alpha - p_v^\alpha\|)}{|F_u^\alpha|}$$

$$d(p_u^\alpha, p_v^\alpha) = \begin{cases} \frac{\lambda}{2 \triangleleft [O_u^\alpha(p_u^\alpha), O_v^\alpha(p_v^\alpha)]} & \|p_u^\alpha - p_v^\alpha\| > \lambda \\ \text{otherwise} & \end{cases}$$

Let us emphasize some properties of the dissimilarity cost. It is useful that it does not depend on the color of the vehicle nor exact location and small changes are allowed. Some examples are shown in Figure 5. Also, one can notice that the score is not inherently symmetric ($d_{uv} \neq d_{vu}$) and the symmetrization will be introduced later.

The dissimilarity cost should not be computed for every two pairs of vehicles as it would be computationally infeasible – there are $N(N - 1)$ combinations. Such computation would be bearable up to thousands or small tens of thousands vehicles, but in our work, we are aiming at millions. Therefore, we use methods of pre-selection to decrease the number of comparisons. We take only first N_B vehicles which will be used for bootstrapping the computation. Hu [38] invariants $(m_1^\alpha, \dots, m_7^\alpha)$ are computed for shapes of these vehicles for each position α . Then, for every pair of vehicles (u, v) from N_B vehicles function ζ , Eq. (2), is evaluated.

$$\zeta(u m, v m) = \sum_{\alpha \in \mathbb{P}} \sum_{i=1}^7 |\operatorname{sgn}_u m_i^\alpha \cdot \log_u m_i^\alpha - \operatorname{sgn}_v m_i^\alpha \cdot \log_v m_i^\alpha| \quad (2)$$

The computation of the function is fast enough to perform $N_B(N_B - 1)$ comparisons even for high N_B . Finally, for each

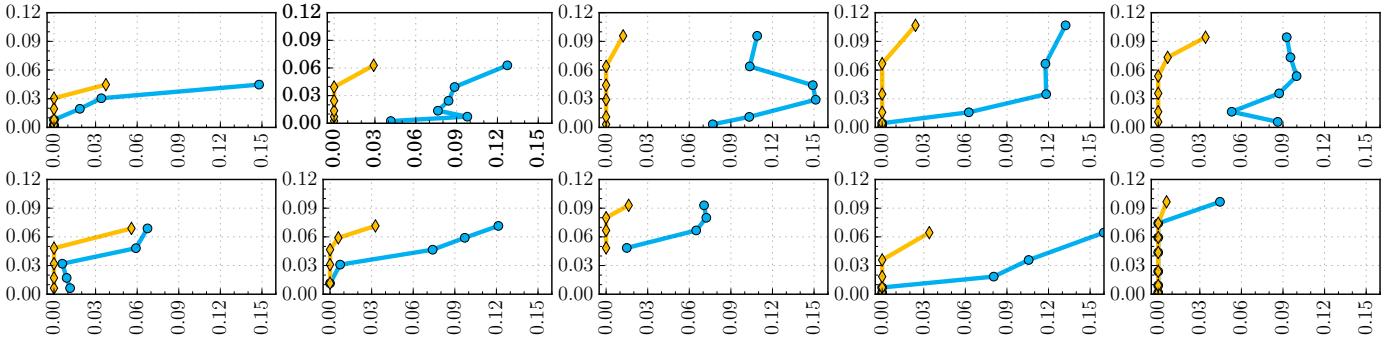


Fig. 7. Grouping algorithm manual evaluation. Each graph represent one camera and lane. The mean false positive rate across 15 largest groups, represented by the X axis, for each settings of thresholds. Y axis represents the fraction of vehicles in the dataset included in clusters with more than 5 vehicles. Yellow curves denote evaluation of make&model grouping. The blue curves represent evaluation when even different submodels of one make&model are considered as different types. The curves are not monotonically increasing as the number of vehicles in the groups may grow and the absolute number of false positives stay the same.

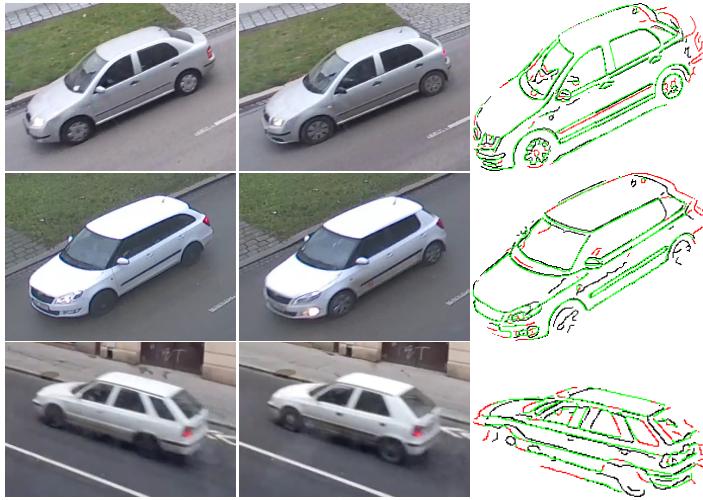


Fig. 8. Rows represent examples of false positives for different groups of vehicles. As it is seen the vehicles are very similar and differ only in the submodel of the vehicle (eg. long version and sedan version).

is processed separately, Fig. 2. The cameras were installed in multiple cities and countries and the types of observed roads range from local ones to large roads with intensive traffic.

The evaluation is focused mainly on automatic grouping and classification accuracy, precision of comparison of vehicles, and accuracy of scene scale inference. The comparison with other methods [1], [2], [3], [4], [5] on fine-grained recognition is not relevant as they usually require the frontal viewpoint.

A. Automatic Grouping Quality

Our first experiment to evaluate the quality of vehicle grouping is based on the quantity of vehicles in the found groups of vehicles and the mean false positive rate across these groups. Figure 7 demonstrates the result where the yellow lines represent the quality of groups. To obtain the evaluation, we performed the grouping on numerous cameras with varying thresholds and for each camera and lane we evaluated the mean false positive rate across the 15 largest (the significant ones) groups and the number of vehicles in the groups of

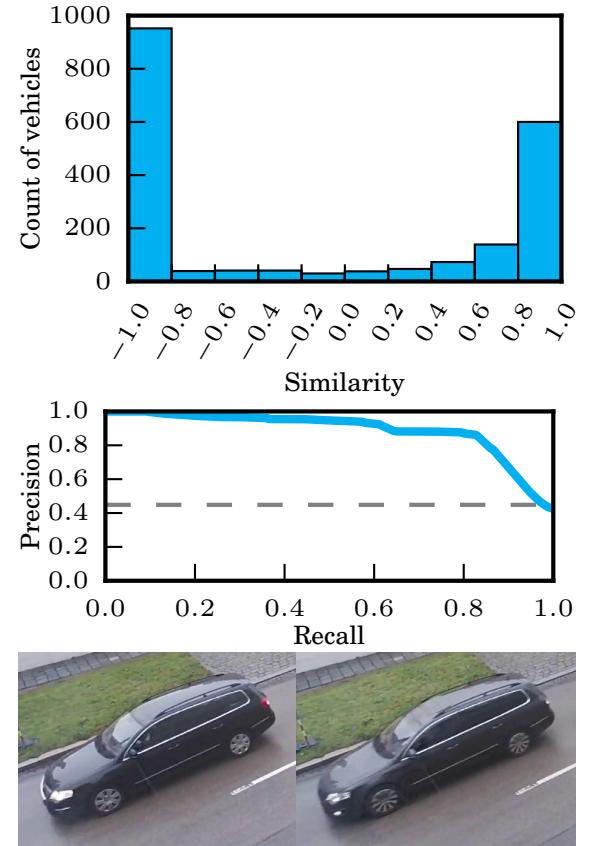


Fig. 9. Evaluation of accuracy of the grouping algorithm on humans annotated pairs of vehicles. People annotated 2,000 pairs of vehicles as “same type” and “different type” while identical types of vehicles should be the same in the make, model, submodel, and model year. **top:** Histogram of normalized annotation responses, where 1.0 means that every annotation for the pair is “same vehicle” and -1.0 denotes the opposite. The middle image shows Precision-Recall Curve for the grouping algorithm when the pairs annotated by human are considered as ground truth, gray dashed line is the baseline. Finally, the bottom image shows a pair of vehicles which are annotated (correctly) as “different type”.

vehicles with more than 5 vehicles. For the experiment, we used $N_B = 20,000$ and $N_H = 10,000$.

In order to examine our method further, we also evaluated

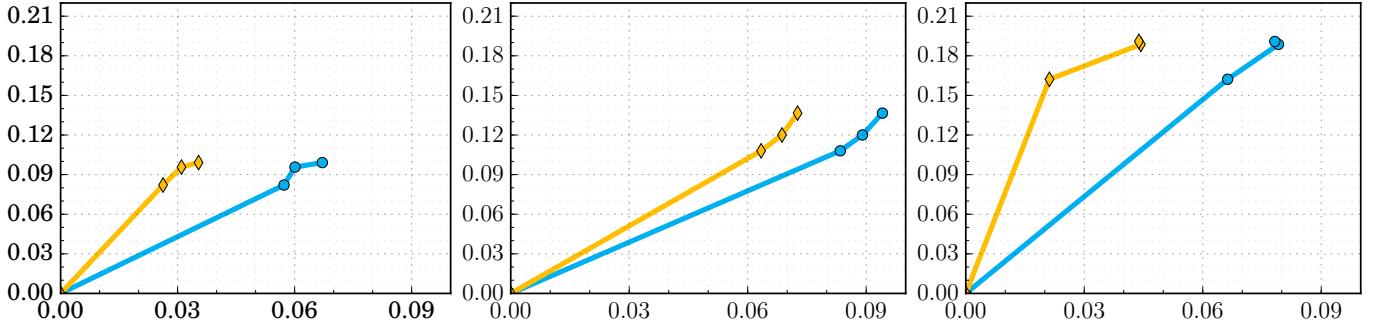


Fig. 10. Evaluation of the classifiers trained on the automatically obtained groups of vehicles: the mean false positive rate across 10 largest groups, represented by the X axis, for each settings of thresholds. The Y axis represents the fraction of vehicles added to the first 10 clusters from the whole dataset. Yellow curves denote evaluation of make&model grouping. Blue curves represent evaluation when even different submodels of one make&model are considered as different types.



Fig. 11. Examples of groups of vehicles. Yellow boxes contain representatives of groups obtained by the proposed unsupervised grouping. Vehicles added by the trained classifiers are in the green boxes. False positives for each group added by the classifiers are in the red boxes. As the images show, our method is robust to color (row 1) and markers variations (row 2), weather (rows 2,3) and lighting changes (rows 2,3) and partial occlusions (see row 4, where vehicles are partially occluded by a tram or by a plastic bag on a tree). Note that the vehicle recognizer works well despite high occlusion, scene clutter and low image quality.

how our algorithm behaves when we consider also different submodels (eg. long version and sedan version) of one make&model as different types. The results are included in Figure 7 as blue curves. This task is really challenging as the submodels are usually very similar and differ only in small subtle details. For examples of such situations see Figure 8.

B. Evaluation of Grouping on a Human-Annotated Dataset

In order to examine how our method will be able to distinguish pairs of identical and different types we used a dataset of 2,000 vehicle pairs annotated by a human. Volunteers were

instructed to recognize same and different types. To be labeled as a match, the pair of vehicles had to match in make, model, submodel, and also model year. The difference of vehicles in different model years can be really small, almost unnoticeable as it can be for example only a small change in shape of lights. One pair was annotated by multiple humans and we collected more than 9 annotations for each pair.

The top graph in Figure 9 shows normalized responses for vehicle pairs while 1.0 denotes that all annotators agreed that the pair contains vehicles of the same type. Value -1.0 denotes that all annotators considered the vehicles' types as

