# Image-based Vehicle Analysis using Deep Neural Network: A Systematic Study

Yiren Zhou, Hossein Nejati, Thanh-Toan Do, Ngai-Man Cheung, Lynette Cheah
Singapore University of Technology and Design
yiren_zhou@mymail.sutd.edu.sg, {hossein_nejati, thanhtoan_do, ngaiman_cheung, lynette}@sutd.edu.sg

*Abstract*—We address the vehicle detection and classification problems using Deep Neural Networks (DNNs) approaches. Here we answer to questions that are specific to our application including how to utilize DNN for vehicle detection, what features are useful for vehicle classification, and how to extend a model trained on a limited size dataset, to the cases of extreme lighting condition. Answering these questions we propose our approach that outperforms state-of-the-art methods, and achieves promising results on image with extreme lighting conditions.

*Index Terms*—Vehicle Classification, Deep Neural Network

## I. INTRODUCTION

Vehicle detection and classification are important parts of Intelligent Transportation Systems. They aid traffic monitoring, counting, and surveillance, which are necessary for tracking the performance of traffic operations. Existing methods use various types of information for vehicle detection and classification, including acoustic signature [1], radar signal [2], frequency signal [3], and image/video representation [4]. The evolution of image processing techniques, together with wide deployment of road cameras, facilitate image-based vehicle detection and classification.

Various approaches to image-based vehicle detection and classification have been proposed recently. Sivaraman and Trivedi [5] use active learning to learn from front part and rear part vehicle images, and achieves 88.5% and 90.2% precision respectively. Chen et al. [6] use a combination of Measurement Based Features (MBF) and intensity pyramid-based HOG (IPHOG) for vehicle classification on front view road images. A rear view vehicle classification approach is proposed by Kafai and Bhanu [7]. They define a feature set including tail light and plate position information, then pass it into hybrid dynamic Bayesian network for classification.

Fewer efforts have been devoted in rear view vehicle classification [7]. Rear view vehicle classification is an important problem as many road cameras capture rear view images. Rear views are also less discriminative and therefore more challenging. Furthermore, it is more challenging for images captured from a distance along multi-lane highways, with possiblity of partial occlusions and motion blur that complicate detection and classification.

We here focus on DNN-based vehicle detection and classification based on rear view images, captured by a static road camera from a distance along a multi-lane highway (Fig. 1). DNN has been applied to many image/video applications [8], [9], [10], [11], [12]. Whilet these methods achieve state-of-the-art on various datasets [13], direct application of them requires a large dataset, that is laborious and expensive to construct. Training of DNN on a small dataset on the other hand, would result in overfitting. Given the difficulty of the original problem (i.e. large in-class variances and ambiguity), reliable modeling based on a small dataset proves even more challanging.

In this work, we propose a combination of approaches to use DNN architectures for this specific problem, building around using the higher layers of a DNN trained on a specific large labeled dataset [14]. There are two approaches to making use of the higher layers of DNN architecture: one is to fine-tune the higher layers of DNN model on our dataset, and the other is to extract the higher layers of DNN architecture as high-level features, and use them for detection and classification. Proposed schemes for our approach are shown in Table I.

Comparing with the state-of-the-art classification methods shows that the vehicle classification methods achieve the highest accuracies. In addition, when coupled with illumination and color transformation and late fusion, the same model retain robustness in classification of poorly lit images *without* fine-tuning or re-training the classification model. Without color transformation, these dark images significantly affect classification results. Our contribution is therefore an approach to train DNN models for vehicle detection and classification on small datasets, and extend their application to cases beyond the content of the original small dataset.

## II. METHODOLOGY

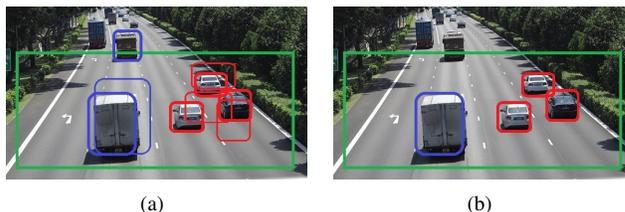### A. Vehicle detection using YOLO



Fig. 1. Examples for vehicle detection approach on a road image. The green rectangle is the selected road region for detection. Red and blue rectangles in (a) are the initial detection results by YOLO model. After remove invalid detection results, the final detection results are shown in (b).

Our dataset images are taken from a static camera along an express way and contain rear views of vehicles on multiple lanes (Fig. 1). We manually label the location bounding boxes for vehicles inside each the road region closest to the camera.

| | Scheme | Dataset [15] |
|---|---|---|
| Vehicle detection | Fine-tune YOLO model | A: 438 road images, |
| | Conventional state-of-the-art | 263 for training, 175 for testing |
| Vehicle classification | Fine-tune Alexnet model | B: 2427 vehicle images, |
| | Alexnet feature extraction | 1440(845 passenger, 595 other) for training, |
| | Conventional state-of-the-art | 987(597 passenger, 390 other) for testing |
| Vehicle classification on dark images | Classification on dark image | C: 257 dark vehicle images, |
| | Classification on transformed image | 223 passenger, 34 other |
| | Late-fusion | |

TABLE I
PROPOSED SCHEMES FOR OUR APPROACH.

A detected vehicle object is valid only if the center of the object is inside the selected road region.

DNN architecture has been widely applied for object detection tasks. Fast R-CNN [8] achieved state-of-the-art result on several datasets. However, a more recent DNN method called YOLO [16] achieved comparable results while being significantly faster. To do the vehicle detection more efficiently, we choose YOLO for our approach.



Fig. 2. Simplified YOLO network structure. A detailed structure can be found in [16].

Fig. 2 shows the simplified structure for YOLO network. The original YOLO network is trained on PASCAL dataset [17] with 20 classes of objects with a probability grid with size *7\*7\*20*. We can increase the size of probability grid to improve detection accuracy. However, it will also increase model complexity, and require higher number of training samples. Here we increase the probability grid to *11\*11*, and to number of classes to 1, resulting in a probability grid with size *11\*11*. Then we fine-tune the last layer of the model with our own road images.

Fig. 1(a) shows that the fine-tuned YOLO model has generates some invalid detection results. We therefore use a post-processing approach to remove these outliers. Outliers are removed based on the following criterion:

$$A \text{ is} \begin{cases} invalid, & \text{if } \exists B \in \text{image}, \ (\frac{Int(A,B)}{Area(A)} > t \| \frac{Int(A,B)}{Area(B)} > t) \\ & \& Conf(A) < Conf(B) \\ invalid, & \text{else if } Center(A) \notin Region(valid) \\ valid, & \text{otherwise} \end{cases}$$

where $A, B$ are two different detected bounding boxes. $Int(A, B)$ is the intersection area of $A, B$, $Area(A)$ is the area of $A$, $t$ is a threshold value, and $Conf(A)$ is the confidence value of $A$ given by YOLO model. $Center(A)$ is the center pixel of $A$, and $Region(valid)$ represents the green rectangle shown in Fig. 1(a). The final detection results after post-processing are shown in Fig. 1(b).

The YOLO model can also perform vehicle classification when we set the number of class to 2, representing passenger and other vehicles. However, the classification accuracy for YOLO is not high enough. We continue to introduce more classification approaches in Section II-B.

## B. Vehicle classification approaches

For vehicle classification, we use the dataset B described in Table I. Vehicle images will be classified into two classes: *passenger* and *other*. Passenger vehicle class includes sedan, SUV, and MPV, other vehicle class includes van, truck, and other types of vehicle. Both classes have large in-class variance. Also the difference between passenger vehicles and other vehicles is not distinctive. These make it difficult to distinguish between these two classes. Fig. 3 shows examples for both vehicle classes. As we can see from the sample images, Fig. 3(a) is MPV, and Fig. 3(b) is taxi. They are both passenger vehicles but different in shape, color, and size. Fig. 3(a) is MPV, and Fig. 3(c) is van. They are in different classes, but similar in shape, color, and size. The classification between passenger vehicles and other vehicles has semantic meanings included, that can only be represented using both low-level and high-level features.
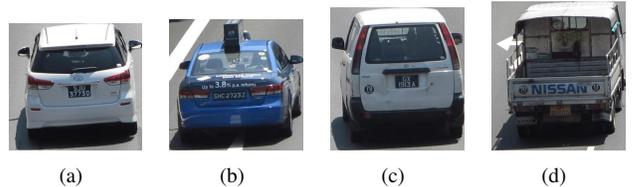


Fig. 3. Vehicle image examples for both classes. (a) passenger. (b) passenger. (c) other. (d) other.

Here we apply two approaches for utilizing DNN architecture: feature extraction, and fine-tuning. For both approaches, we adopt Alexnet [9] model as DNN architecture.

For each vehicle image detected from Section II-A, we resize it to $256 \times 256$, make it valid Alexnet input. Then the resized image is passed into Alexnet. Fig. 4 shows the structure of Alexnet. Alexnet has 5 convolutional layers (named as conv1 to conv5) and 3 fully-connected layers (named as fc6, fc7, fc8). Each convolutional layer contains multiple kernels, and each kernel represents a 3-D filter connected to the outputs of the previous layer. For fully-connected layers, each layer contains multiple neurons. Each neuron contains a positive value, and it is connected to all the neurons in previous layer.
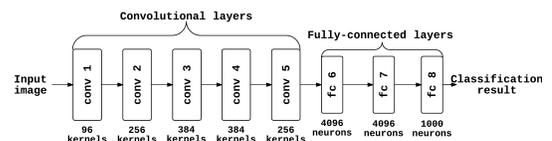


Fig. 4. Structure of Alexnet.

*1) Feature extraction using Alexnet:* Here we extract the third last and second last fully connected layer (i.e. layer fc6 and fc7) in Alexnet as the generic image representation

(to be justified later). Each image representation is a 4096-dimension vector, obtained from the 4096 neurons in layer fc6 (or fc7). Here we consider the extracted layer as a feature vector $f = [f_1, f_2, ..., f_{4096}]$. After we obtain the deep feature vector, SVM with linear kernel is used for classification.

Different layers in a Deep Neural Network (DNN) are often considered to have different level of features. The first few layers contain general features that resemble Gabor filters or blob features. The higher layers contain specific features, each representing a particular class in dataset [14]. Thus features in higher layers are considered to have higher level vision information compared to general features in base layers. To understand this in our particular problem, Fig. 5 shows several average images we obtained from vehicle images. Given a specific feature $f_i$ we extracted from Alexnet, we sort all the vehicle images based on value of $f_i$. The images that have highest values on this feature are chosen. Then we calculate the average image of these images. The 4 images in Fig. 5 represents average images for 4 different features (i.e. $f_{i_1}, ..., f_{i_4}$, here $i_1, .., i_4 \in \{1, ..., 4096\}$). We can recognize specific types of vehicles from these average images. Fig. 5(a) represents a specific type of normal sedan. Fig. 5(b) is taxi. Fig. 5(c) is van. And Fig. 5(d) represents truck. Human can easily associate these average images to certain types of vehicles, meaning that the features related to these images contain high-level visualization information related to semantic meanings of each class.
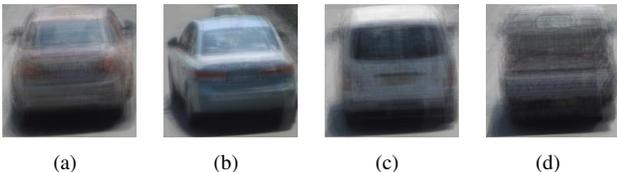


(a)　　　　(b)　　　　(c)　　　　(d)

Fig. 5.　Average image of the vehicles with high values on a specific feature.

*2) Fine-tuning Alexnet on our dataset:* Another approach to make use of the high-level information in DNN is to fine-tune the DNN model on our dataset. Alexnet [9] is trained with 1000 classes. To match our dataset with 2 classes, we change the size of fc8 layer of Alexnet from 1000 to 2. Then we use Alexnet model trained on ILSVRC 2012 dataset to fine-tune on our dataset. In order to prevent overfitting, the parameters from layer conv1 to layer fc6 is fixed. After the fine-tuning, the model is tested on testing set with 987 images.

*C. Vehicle classification on dark images*

There is also a need for vehicle classification on dark images. Fig. 6(a) shows an image taken during the night, where classification is more challenging due to poor lighting. One approach to improve accuracy for dark vehicle image classification is to train the model on dark images, however, it is not feasible when we have a limited number of dark image samples. Here we propose a method to use model trained on normal images to classify dark images.

Dark image often comes with low contrast, color displacement, and high noise that would significantly affect image quality. By applying color transformation method we
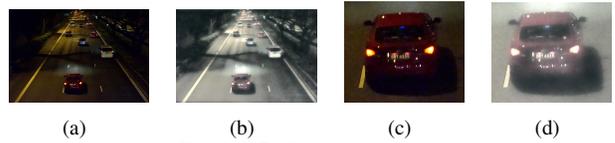


(a)　　　　(b)　　　　(c)　　　　(d)

Fig. 6.　Dark image examples.

can increase image contrast and fix color displacement. In addition, the classification model is trained on *normal* images, so we want to have a transformed image close to *normal* images. Here we make use of a high-level scene transformation method [18] to transform *night* image back to *normal* image. This transformation model are trained from scene image dataset with different lighting, weather, and seasons. It can conduct high-level scene transformation, including the transformation between different lighting conditions.

Fig. 6(b) is the transformed image by using [18]. If we see details in passenger vehicle Fig. 6(c), the contrast level is low, and there also exist color displacement and noise. By using scene transformation, the contrast level of Fig. 6(d) increases, color displacement and noise still remains. However, from the transformed image Fig. 6(d) we can observe that the color displacement is different from Fig. 6(c). From these results we propose a late fusion method to utilize both the dark image and the transformed image for classification.

The late fusion method is listed as follows:

$$Label_{fused} = \arg\max_i Conf(i, j), i \in \{passenger, other\}$$
$$j \in \{original, transformed\}$$
(1)

where $Conf(i, j)$ is the confidence score of class $i$ from image $j$, generated from SVM model. For Alexnet fine-tuning model, the confidence score is generated by the softmax layer. The idea is to select the classification result with highest confidence on both night image and transformed image.

We test the scene transformation and late fusion methods on a night image dataset C in Table I. The test results are reported in Section III-B2.

## III. Experimental results

In this section the experimental results of the proposed vehicle detection and vehicle classification method are presented. The size of road images in our dataset is $4184 \times 3108$. For the vehicle detection process, the road images are resized to $448 \times 333$ using Bicubic interpolation. After vehicle detection, we map the vehicle regions back to original road images, and crop vehicle image in original resolution. Typical resolution of vehicle images is around $500 \times 500$. For vehicle classification, all vehicle images are resized to $256 \times 256$ to pass into Alexnet, for other embedding methods, we use vehicle images with original resolution. The vehicle detection method is implemented in darknet [19]. For vehicle classification method, the feature extraction and fine-tuning of Alexnet is under Caffe framework [20], other feature embedding methods and SVM are implemented in MATLAB.

*A. Vehicle detection experiment*

We train and test the YOLO detection model on dataset A described in Table I. Here we compare the result with

| Detection result | | YOLO fine-tune | | DPM [21] | |
|---|---|---|---|---|---|
| | | positive | negative | positive | negative |
| Ground truth | positive | 921 | 185 | 932 | 149 |
| | negative | 66 | - | 55 | - |

TABLE II
COMPARISON OF DETECTION RESULT.

another state-of-the-art detection method DPM [22], [21]. Among 987 testing images, 921 are successfully detected, 185 detected images are invalid images. Vehicle detection precision is 93.3%, and recall is 83.3%, compared with 94.4% and 86.2% by DPM method.

### B. Vehicle classification experiment

| Accuracy (%) | Cars vs vans | Sedans vs taxis | Sedans vs vans vs taxis |
|---|---|---|---|
| PCA+DFVS [23] | 98.50 | **97.57** | 95.85 |
| PCA+DIVS [23] | 99.25 | 89.69 | 94.15 |
| PCA+DFVS+DIVS [23] | - | - | 96.42 |
| Constellation model [24] | 98.50 | 95.86 | |
| Alexnet-fc6-SVM | **99.50** | 97.27 | **97.36** |
| Alexnet-fc7-SVM | 99.25 | 96.67 | 94.75 |

TABLE III
ACCURACY COMPARISON ON PUBLIC DATASET. REPORTED RESULTS
FROM [23] ARE USED.

*1) Experiment on public dataset:* To compare our approach with other classification methods, we perform our approach on a public dataset provided in [24]. We use same experiment setting in [23] to perform fair comparison. There are three types of vehicles in this dataset: sedans, vans, and taxis. Following [23], three experiments are performed: *cars* vs *vans*, *sedans* vs *taxis*, and *sedans* vs *vans* vs *taxis*. Note that sedans and taxis are all regarded as cars.

We apply feature extraction on Alexnet model for classification. From each vehicle images, we extract two feature vectors (from layer fc6 and fc7) with 4096 dimensions using Alexnet. Then, linear-SVM with is applied for classification.

Table III shows accuracy comparison among Alexnet-based methods and other state-of-the-art methods. Alexnet-fc6 feature achieves best accuracy on *cars* vs *vans*, and *sedans* vs *vans* vs *taxis* classification, and second-best accuracy on *sedans* vs *taxis* classification. These results show the effectiveness of Alexnet features on vehicle classification problem.

*2) Experiment on our dataset:* The vehicle classification approaches are trained and tested on dataset B[1] in Table I. All results are calculated using class-balanced accuracy as shown below:

$$Acc_{bal} = \frac{\frac{Correct(pass)}{Size(pass)} + \frac{Correct(other)}{Size(other)}}{2} \qquad (2)$$

where $Correct(pass)$ is the number of correct prediction in passenger class, and $Size(pass)$ is the total number of images in passenger class.

Here we compare the performance of the DNN with state-of-the-art image description methods: Fisher vector [25], FAemb [26], [27], and Temb [28] with SIFT descriptor[2].

[1]The dataset can be downloaded via link [15].

[2]Unable to run the code from [23], we did not include their methods.

From each vehicle image, we extract a feature vector (fc6 or fc7) with 4096 dimensions using Alexnet. Another alternative method is to concatenate fc6 and fc7 to get 8192 dimensions feature vector. For Alexnet fine-tuning, we directly use the fine-tuned model to classify vehicle images.

For comparison with state-of-the-art methods, from each vehicle image, we first compute SIFT descriptors (each having 128 dimensions) of the image. Then different embedding features are generated based on SIFT. Generated fisher vectors have 4k or 8k dimensions. Temb and Faemb have around 8k dimensions.

For both Alexnet extracted features and other methods, we use linear-SVM to train the classifier. The SVM is trained on dataset B in Table I. The trained model is also tested on a dark image dataset C.

| Accuracy (%) | Dims | Normal images | Dark images | Transformed images | Late fusion |
|---|---|---|---|---|---|
| Fisher-vec-4k [25] | 4096 | 93.55 | 57.3 | 60.06 | 60.85 |
| Fisher-vec-8k [25] | 8192 | 93.3 | 58.98 | 62.37 | 61.87 |
| FAemb [26], [27] | 8280 | 89.76 | 69.7 | 60.13 | 68.67 |
| Temb [28] | 8192 | 87.81 | 65.35 | 56.81 | 64.55 |
| Alexnet-fc6 | 4096 | 96.95 | **74.68** | **79.32** | **85.41** |
| Alexnet-fc7 | 4096 | 96.44 | 57.93 | 64.91 | 67.51 |
| Alexnet-fc6&7 | 8192 | **97.35** | 73.65 | 77.97 | 84.16 |
| Alexnet-fine-tune | - | 97.15 | 52.76 | 52.08 | 52.52 |

TABLE IV
ACCURACY COMPARISON ON OUR DATASET.

Table IV shows the accuracy comparison for all methods. The best result is achieved by concatenating fc6 and fc7 layers of Alexnet. Fine-tuning also achieves good result for classification on normal images.

For dark image results, all methods have suffered from severe accuracy degradation. fc6 feature has achieved best result on dark images. We can see the fc7 feature has much lower result compared to fc6, indicating that fc6 is more robust to low contrast and color displacement. It is also interesting to see that fine-tune Alexnet model has poor result on dark images. The fc7 model and fine-tuned Alexnet model is fitted into normal image classification, and cannot generalize to dark image classification.

All Alexnet features have improvement on transformed images, indicating the effectiveness of scene transformation on dark images for classification. Other state-of-the-art feature embedding methods are not benefited from the transformation, because these features are generated from SIFT feature, and SIFT feature does not have strong relationship with color transformation of the image.

The late fusion results are shown in the last column. We can see that the Alexnet features have improvements after we conduct late fusion on dark and transformed images. The best result is achieved by fc6 feature for about 85%.

## IV. CONCLUSION

We have investigated DNN approaches for both vehicle detection and classification using a limited size dataset. For detection, we fine-tune a DNN detection model for vehicle detection and achieved good result. For classification, we

evaluate both fine-tuning and feature-extraction method, the result outperformed state-of-the-art.

We further proposed methods to use scene transformation and late fusion techniques for classification on poor lighting conditions, and achieved promising results without changing the classification model. Our approach is therefore have the potential to be used for training on limited size datasets and be extended to different cases such as various lighting conditions.

## REFERENCES

[1] Kangyan Wang, Rui Wang, Yutian Feng, Haiyan Zhang, Qunfeng Huang, Yanliang Jin, and Youzheng Zhang, "Vehicle recognition in acoustic sensor networks via sparse representation," in *IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*. IEEE, 2014, pp. 1–4.

[2] Heong-tae Kim and Bongsob Song, "Vehicle recognition based on radar and vision sensor fusion for automatic emergency braking," in *13th International Conference on Control, Automation and Systems (ICCAS)*. IEEE, 2013, pp. 1342–1346.

[3] Troy R McKay, Carl Salvaggio, Jason W Faulring, Philip S Salvaggio, Donald M McKeown, Alfred J Garrett, David H Coleman, and Larry D Koffman, "Passive detection of vehicle loading," in *IS&T/SPIE Electronic Imaging*. International Society for Optics and Photonics, 2012, pp. 830511–830511.

[4] Pradeep Kumar Mishra and Biplab Banerjee, "Multiple kernel based knn classifiers for vehicle classification," *International Journal of Computer Applications*, vol. 71, no. 6, 2013.

[5] Sayanan Sivaraman and Mohan M Trivedi, "Real-time vehicle detection using parts at intersections," in *15th International IEEE Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2012, pp. 1519–1524.

[6] Zezhi Chen, Tim Ellis, and SA Velastin, "Vehicle detection, tracking and classification in urban traffic," in *15th International IEEE Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2012, pp. 951–956.

[7] Mehran Kafai and Bir Bhanu, "Dynamic bayesian networks for vehicle classification in video," *IEEE Transactions on Industrial Informatics*, vol. 8, no. 1, pp. 100–109, 2012.

[8] Ross Girshick, "Fast r-cnn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.

[9] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[10] S. Song, V. Chandrasekhar, B. Mandal, L. Li, J. H. Lim, G. S. Babu, P. P. San, and Ngai-Man Cheung, "Multimodal multi-stream deep learning for egocentric activity recognition," in *Proc. of IEEE CVPR - 4th Workshop on Egocentric (First-Person) Vision*. IEEE, 2016.

[11] Thanh-Toan Do, Anh-Dzung Doan, and Ngai-Man Cheung, "Learning to hash with binary deep neural network," in *Proc. of European Conference on Computer Vision (ECCV)*. IEEE, 2016.

[12] V. Pomponiu, H. Nejati, and N.-M. Cheung, "Deepmole: Deep neural networks for skin mole lesion classification," in *Proc. of ICIP*. IEEE, 2016.

[13] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson, "Cnn features off-the-shelf: an astounding baseline for recognition," *arXiv preprint arXiv:1403.6382*, 2014.

[14] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson, "How transferable are features in deep neural networks?," in *Advances in Neural Information Processing Systems*, 2014, pp. 3320–3328.

[15] "Vehicle dataset," https://goo.gl/Jpu2Ox.

[16] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi, "You only look once: Unified, real-time object detection," *arXiv preprint arXiv:1506.02640*, 2015.

[17] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results," http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html.

[18] Pierre-Yves Laffont, Zhile Ren, Xiaofeng Tao, Chao Qian, and James Hays, "Transient attributes for high-level understanding and editing of outdoor scenes," *ACM Transactions on Graphics (proceedings of SIGGRAPH)*, vol. 33, no. 4, 2014.

[19] Joseph Redmon, "Darknet: Open source neural networks in c," http://pjreddie.com/darknet/, 2013-2016.

[20] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2014, pp. 675–678.

[21] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part based models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010.

[22] R. B. Girshick, P. F. Felzenszwalb, and D. McAllester, "Discriminatively trained deformable part models, release 5," http://people.cs.uchicago.edu/ rbg/latent-release5/.

[23] Amol Ambardekar, Mircea Nicolescu, George Bebis, and Monica Nicolescu, "Vehicle classification framework: a comparative study," *EURASIP Journal on Image and Video Processing*, vol. 2014, no. 1, pp. 29, 2014.

[24] Xiaoxu Ma and W Eric L Grimson, "Edge-based rich representation for vehicle classification," in *Tenth IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2005, vol. 2, pp. 1185–1192.

[25] Jorge Sánchez, Florent Perronnin, Thomas Mensink, and Jakob Verbeek, "Image classification with the fisher vector: Theory and practice," *International journal of computer vision*, vol. 105, no. 3, pp. 222–245, 2013.

[26] Thanh-Toan Do, Quang D Tran, and Ngai-Man Cheung, "Faemb: a function approximation-based embedding method for image retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3556–3564.

[27] Thanh-Toan Do and Ngai-Man Cheung, "Embedding based on function approximation for large scale image search," *CoRR*, May 2016.

[28] Hervé Jégou and Andrew Zisserman, "Triangulation embedding and democratic aggregation for image search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3310–3317.