

VPDS: An AI-Based Automated Vehicle Occupancy and Violation Detection System

Abhinav Kumar,^{1*†} Aishwarya Gupta,^{2*} Bishal Santra,^{3*†} Lalitha KS,^{2*}

Manasa Kolla,^{2*} Mayank Gupta,^{2*} Rishabh Singh^{2*}

¹School of Computing, University of Utah, ²Computer Vision and Media Analytics Group, Conduent Labs

³Indian Institute of Technology Kharagpur

abhinav.kumar@utah.edu, aishwarya.gupta@conduent.com, bsantraigi@gmail.com, lalitha.ks@conduent.com
manasa.manasa@conduent.com, mgmayank18@gmail.com, rishabh.singh@conduent.com

Abstract

High Occupancy Vehicle/High Occupancy Tolling (HOV/HOT) lanes are operated based on voluntary HOV declarations by drivers. A majority of these declarations are wrong to leverage faster HOV lane speeds illegally. It is a herculean task to manually regulate HOV lanes and identify these violators. Therefore, an automated way of counting the number of people in a car is prudent for fair tolling and for violator detection.

In this paper, we propose a Vehicle Passenger Detection System (VPDS) which works by capturing images through Near Infrared (NIR) cameras on the toll lanes and processing them using deep Convolutional Neural Networks (CNN) models. Our system has been deployed in 3 cities over a span of two years and has served roughly 30 million vehicles with an accuracy of 97% which is a remarkable improvement over manual review which is 37% accurate. Our system can generate an accurate report of HOV lane usage which helps policy makers pave the way towards de-congestion.

Introduction

Intelligent Transportation Systems (ITS) improve safety and mobility through the integration of sensing, computational power and advanced communications into the transportation infrastructure (Xu et al. 2014). Such systems enable efficient management of lanes by incorporating various aspects like carpooling, tolling, traffic management and transit in a multi-purpose roadway. This creates novel avenues for agencies in terms of congestion pricing in order to generate revenue and manage demand dynamically. Population growth has resulted in heavy congestion on highway lanes, leading to both economic and environmental concerns. Recent statistics reveal a monotonic increase in the number of vehicles on highways from 193 million in 1990 to approximately 268.8 million vehicles registered in 2016 in USA (sta 2018).

HOV lanes are standard car-pool lanes where a minimum of two (HOV2+) or three (HOV3+) vehicle occupants are required to use the lane legally, making them lesser congested and enabling efficient rapid transit (Daley et al. 2011). In order to avoid congestion yet still encourage car-pooling,

agencies allow cars with a single occupant to use carpool lanes by paying a toll. These High Occupancy Toll (HOT) lanes have a variable toll. Generally, agencies try to maintain a minimum speed of 45 miles per hour (mph) on such lanes, so they need to monitor the number of single occupant vehicles allowed to enter HOV lanes. This is done typically by charging a toll adjusted to the dynamic congestion in the lane. The toll price is increased if the average traffic speed in the HOT lane decreases below the accepted minimum speed.

However, to gain the benefits offered by HOV/HOT lanes, the entry rules (number of occupants in the vehicle) need to be enforced vigilantly. The declaration of the occupancy status of a vehicle is a voluntarily compliance by the driver. But in most cases, this declaration is falsified in order to avoid the toll. The current practice is to rely on visual inspection by road-side officers to enforce these rules, but the process is found to be inefficient, costly and potentially dangerous (Artan et al. 2016). Typical violation rates can exceed 50-80%, while manual enforcement rates are typically less than 10% (Schijns and Mathews 2005). While tagging genuine passengers as violators causes discomfort among the consumer base, allowing too many violators in an HOV lane nullifies its purpose. In either case, the loss is incurred by the transportation agency providing the service. Therefore, an automated way of counting the number of occupants of a vehicle is extremely necessary for fair tolling and violator detection.

In this paper, we address the above challenge by proposing a Vehicle Passenger Detection System (VPDS) - a deep neural network based solution for counting the number of passengers inside a vehicle by processing its front and side images. Our contribution to the problem is two-fold. We first apply a state-of-the-art object detection technique, YOLOv3 (Redmon and Farhadi 2018), for the front and the rear window detection. This enables a fast and accurate localization of the region of interest (ROI) in the image. Next, we perform classification (to count number of people) separately for the front and rear regions of interest, i.e. windows using GoogleNet (Szegedy et al. 2015) based models. Further, we combine the outputs to arrive at the decision of HOV/not-HOV violator for a vehicle. We show the robustness of the proposed solution in terms of better performance to existing approaches within the constraints of poor image quality and significant external factors like illumination, occlusion and traffic congestion.

*Equal contribution by all the authors

†This work was carried out at Conduent Labs

Copyright © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

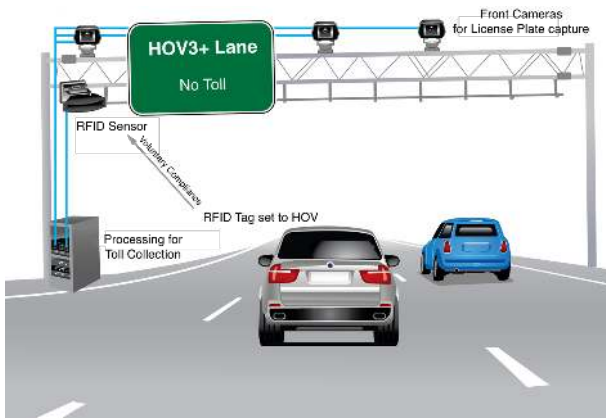


Figure 1: Present: Counting done with an RFID pass.

Application Description

The setup widely used in transportation for HOV lane enforcement is shown in Figure 1. Users of HOV lanes toggle an RFID tag based while entering the tolling booth and the system relies on self compliance by the driver. However, statistics reveal that 80% of the times the driver gives a wrong declaration to avoid toll (Schijns and Mathews 2005). Therefore, automated methods to identify and fine the violators have to be developed which will result in better compliance in the usage of car pooling lanes.

Computer Vision (CV) through Artificial Intelligence (AI) is the most effective way of developing an automated vehicle occupancy counting system. Figure 2 shows the entire pipeline for identifying HOV violators. First, the front and the rear seat¹ images are captured by two cameras, one of which is aimed at the oncoming traffic, while the other is set perpendicular to it. These two images are then processed using AI based approaches. In case the system finds a mismatch between the number of passengers in the vehicle declared during voluntary compliance and the number detected by the system, the driver is fined with the help of the license plate information of the vehicle. License plate recognition is a separate problem which has been tackled in (Bulan et al. 2017). In this paper, the problem addressed is to count the number of people seated in a vehicle and classify the vehicle as an HOV3+ violator or a non-violator.

Background

The images captured by the hardware contains information from the entire scene. The front and side windows of a vehicle form a limited region of the image only. As such, we need to crop the relevant portion of the images which is helpful towards the task of counting passengers. Thus, the first step in the framework is extraction of ROI.

Previous works (Xu et al. 2014) in occupancy detection use DPM (Felzenszwalb et al. 2010) for ROI extraction, but there are sizable limitations to it. DPM relies heavily on image preparation, cropping and in general the saliency of the

¹The terms “rear” (rear seats) and “side” (side view of rear seats) have been used interchangeably.

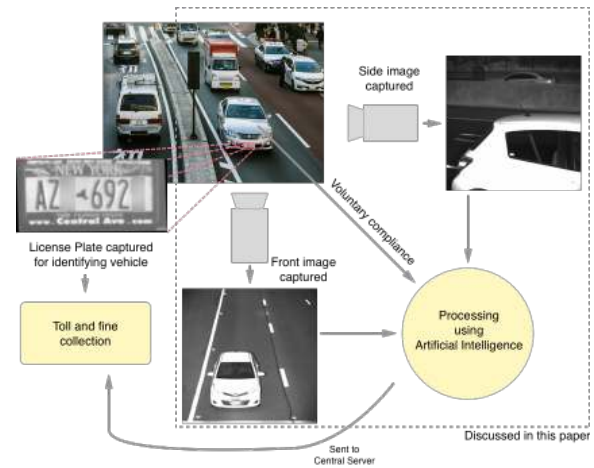


Figure 2: Proposed vehicle occupancy count processing pipeline. This involves capturing and localizing the front and the rear images by two distinct cameras and processing them using AI.

vehicle in the pre-processed image. Therefore, in this paper, we propose ROI extraction using YOLOv3 (Redmon and Farhadi 2018), which does not have such heavy dependencies on image pre-processing.

The next step in the pipeline is to count the number of people in vehicle using these ROIs. Occupancy Detection methods can be broadly classified into three main categories viz, detection, feature and density based methods as explained below.

Detection-based methods Face detection has been extensively explored for counting people by detecting passenger faces using a pixel threshold (Wang, Xu, and Paul 2015). Skin detection (Hao, Chen, and Li 2006) can process front image using a color skin model to coarsely detect the facial region. Occupancy can also be estimated by detecting the empty seats in a vehicle (Fan et al. 2013). However, passengers in a vehicle do exhibit arbitrary poses. Especially in side view images, the visibility of faces is a prominent issue given frequent occlusion. In such cases, detection based methods do not yield convincing results.

Feature-based methods These methods synthesize features which capture the difference between a passenger and his/her surroundings. Some approaches have explored distance-based metrics between descriptors in order to discriminate between images having only the driver or both the driver and the passenger in the front image of a vehicle (Xu, Paul, and Perronnin 2017). Some works have shown superior performance of classification with Fisher Vectors (Perronnin and Dance 2007), (Perronnin, Sánchez, and Mensink 2010) to DPM based models (Artan and Paul 2013). These features, however, fail to capture the variability in low resolution/lesser informative images, which is even more prominent in a real setting.

Density-based methods These methods aim to estimate the count of people in an image by learning to create a density

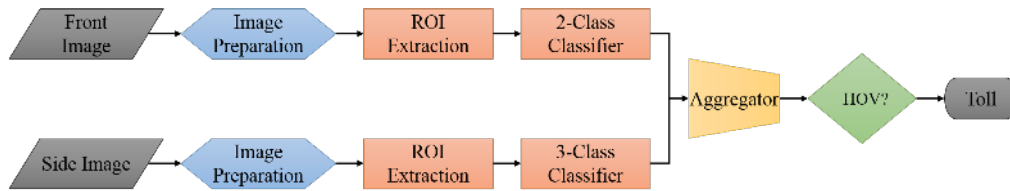


Figure 3: Flowchart of passenger counting and HOV violator detection using Computer Vision techniques. ROI extraction from front and rear images is followed by individual classifiers, results of which are aggregated for the final decision.

map of the input image (Sindagi and Patel 2018). This entire spectrum of approaches fail to estimate people count in low density scenarios.

Parts of the proposed framework were earlier reported in (Xu et al. 2014), (Artan et al. 2016) and (Wshah et al. 2016). This paper claims sufficient novelty and improvements over these previous works. (Xu et al. 2014) uses DPM for localization and Support Vector Machines (SVM) for classification while (Artan et al. 2016) uses DPM for windshields extraction from images and feature-specific models for classification. (Wshah et al. 2016) proposes a solution for detecting whether there is a passenger present in a vehicle or not, which is essentially a binary classification and thus can only be used for figuring HOV2+ violations. Our work not only improves detection, but also counts the exact number of people in a vehicle by incorporating multi-class classification in the rear images. This makes our solution suitable for detecting HOV3+ violations as well.

Usage of AI Technology

Figure 3 describes the vehicle occupancy counting system formulated as a CV problem. The front and rear images are passed through two separate processing streams whose results are combined to obtain the HOV decision. Both these streams have ROI extraction and classification (human counting) modules in common. Each of these modules are explained in the following subsections.

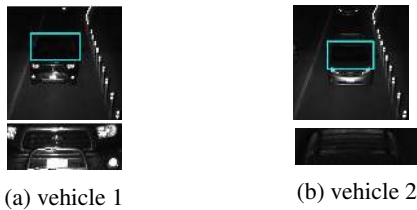


Figure 4: Comparison of ROI detection by YOLOv3 (top) and DPM (bottom). Clearly, the DPM based method fails.

ROI Detection

As previously mentioned, ROI extraction is a necessary step before classification. We evaluate the performance of DPM and YOLOv3 on images from a particular day without any image pre-processing (crop, rotation or scaling). We found that DPM was able to find the correct ROIs in 52% of the

cases, while YOLOv3 was correct in more than 96% cases. The correctness of ROIs was determined based on an Intersection Over Union (IOU) threshold with the ground truth labels. Further, forward pass on YOLOv3 is 10 times faster than on DPM. Figure 4 depicts that YOLOv3 detects the windshields accurately despite the presence of other similar entities like sunroofs, which in the case of DPM have to be removed manually by fine-tuning preprocessing parameters for each camera. This was one of the main drawbacks of the solution proposed in (Wshah et al. 2016). Even after manual parameter tuning, the detection accuracy of DPM based models is lesser than YOLOv3. Thus, YOLOv3 is better for real-time applications like transportation in terms of both time and accuracy.

Person Counting

The passengers in a vehicle can be divided into front and rear seat occupants. ROIs extracted for each view are processed by two separate CNNs. The front row can contain only one person apart from the driver, thus the problem turns into a binary classification of presence/absence of a passenger along with the driver. We use GoogleNet in conjunction with YOLOv3 for this task which gives an accuracy of more than 97%. Our method is a significant improvement over the DPM-enabled classifier which has an accuracy of 94.6% in the best case.

The side image classification is a more challenging task because of several reasons. First, it is a four-class classification problem since the number of passengers can range from zero to three. However, for HOV3+ predictions, three classes (zero, one and more than one classes) would suffice. Second, a heavily skewed dataset results in a bias in predictions. The rear seats are empty in 66.63% of cases, has one occupant 22% of the times, while it has more than one occupant in only 11.37% cases. Finally, since the images are from the side, the passengers are quite susceptible to occlusion by one another. These challenges result in errors and demand an AI-engine robust to such externalities. Thus, we evaluated three popular CNN architectures viz. GoogleNet (Szegedy et al. 2015), ResNet (He et al. 2016) and VGGNet (Simonyan and Zisserman 2014) after oversampling class 1 (1 passenger) and class 2 (2 or more passengers) to match the number of samples in class 0 (no passenger). The front and rear counts obtained from their respective deep CNNs are added to get an estimate of occupancy of a vehicle.

Application Use and Payoff

According to the global traffic survey (inr 2017), drivers in New York spent 91 peak hours stuck in traffic. This traffic congestion will cost an average \$100 billion over the next five years. HOV lanes minimize the delay caused due to traffic congestion by promoting car pooling and hence reduce the number of cars on the highway. These lanes reduce the average travel time by 70% for vehicles moving on HOV lanes. Additionally, for vehicles on the usual lanes the delay reduces by 50%. Further, the cars which have less number of people can still take HOV lanes during congestion by paying a toll. In fact, 33% of vehicles choose to use HOV lanes. The automated vehicle occupancy counting system as described in this paper proves to be extremely effective for non-intervening functioning for huge volumes of traffic flowing across the highways throughout the day all round the year with negligible maintenance.

<i>Model</i>	<i>Precision</i>	<i>Recall</i>	<i>Accuracy</i>	<i>Group Accuracy</i>
Deep Classification Models				
ResNet-50	83.27	77.27	92.57	84.73
VGG-16	88.58	85.29	95.01	88.32
GoogleNet	89.01	83.66	94.82	88.68
Person Detection Models				
YOLOv3	87.20	77.93	93.50	86.20

Table 1: Performance comparison (in percentage) of different counting models - HOV3+ and group . (front + rear)

Table 1 shows the results of using different CNNs on the person counting task. In addition to the standard precision, recall and accuracy, we also report the *group accuracy*. Group accuracy is calculated when individually both the front and rear passengers are counted correctly by the nets rather than the total count of HOV3+ violators or non-violators.

The system's performance over the day is fairly constant and in the range of 94% to 96% accuracy. There is a relative drop in performance during noon and early-night hours. The former can be attributed to significant glare present due to the sun which hinders the visibility of human faces, especially in the front window. The latter variation can be safely ignored as the results are within margin of error and the traffic volume falls considerably post 9 PM.

We evaluated the time based performance of the three models as well. ResNet-50 takes 55 ms per image, VGG-16 takes 26 ms per image and GoogleNet takes 15 ms per image on a GPU. Our proposed system based on GoogleNet takes much less time per image than VGG-16. This in conjunction with the high accuracy achieved (Table 1) demonstrate the effectiveness of the overall system. The numerical figures corroborate our claim of fast processing without compromising on accuracy of decision. Thus the proposed system is efficient for high traffic flow.

The variation in system level accuracy and yield in terms of the system confidence threshold can be observed in Figure 5. The system level confidence is formulated using the individual front and side GoogleNet classifier confidences.

This trend facilitates the use of decision threshold as a design parameter to obtain a better violator detection rate in terms of accuracy. This also caters to the transportation agency to have control over the false positive rate (non-violators detected as violators) of the system. The increasing trend of accuracy with respect to the confidence threshold also illustrates the correctness of the classifier confidence. On the other hand the yield curve acts as another system performance measure showcasing the fraction of vehicles classified for varying confidence thresholds.

Table 2 shows the precision, recall and accuracy achieved by three classification CNNs and YOLOv3 as a person detection system on the side images. The overall accuracy of rear classification is highest for GoogleNet. The count in the case of rear images is less accurate compared to the front images since the CNN gets confused between the images with one versus more than one person sitting in the back due to occlusion. This is confirmed by the visualization of the features learned by GoogleNet for the rear classification task which is shown in Figure 5d. The blue cluster corresponding to empty seats is well-separated from the highly overlapping red and green clusters. This depicts the high accuracy in detecting the presence/absence of a person by the AI model. However, it gets confused in demarcation between one or more passengers owing to significant occlusion when seen from side.

<i>Model</i>	<i>Precision</i>	<i>Recall</i>	<i>Accuracy</i>
Deep Classification Models			
ResNet-50	72.21	71.46	86.87
VGG-16	81.22	79.35	91.04
GoogleNet	81.60	79.32	91.08
Person Detection Models			
YOLOv3	79.92	72.77	88.40

Table 2: Performance comparison of the models for the rear passenger count (in percentage)

In addition, we also plot the yield versus accuracy curve for the rear classification models in Figure 5c. It is consistent with the table 2 that GoogleNet is more accurate in counting people than VGGNet and ResNet for the same yield.

The proposed system is very robust to certain factors that we showcase. The rear seat looks empty in Figure 6a, but histogram equalization shows a baby on a safety seat in Figure 6b. The baby passenger, although not visible to human eyes, is correctly identified as one person by the system.

The performance of the system is slightly affected in the presence of infants in the rear seat, as they are often occluded by the safety carrier. Given the ambiguity in the presence of a child, as shown in Figure 6c, we labelled such data containing only a safety carrier as "No occupant". However, whenever the child is partially visible, the system correctly recognizes him/her as a person as in Figure 6d. Various other scenarios shown in Figure 7 like a cat present at the back (0 person), occlusion of face, only hair of a passenger being visible or pose variation which lead to failure of manual counting process, are efficiently handled by VPDS with high accuracy.

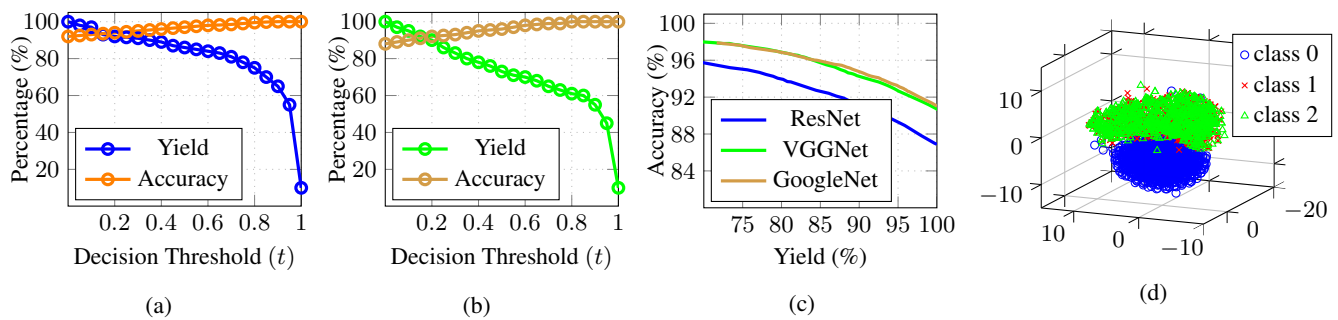


Figure 5: Variation of accuracy and yield of VPDS with threshold on Confidence Score for (a) HOV3+ and (b) Non-HOV3+ predictions. The final confidence score is obtained by multiplying the softmax output of front and the rear classifier. (c) Yield vs Accuracy curve for the rear passenger counting using three popular CNN architectures. (d) A 3-D t-SNE plot of GoogleNet features for the rear classification task. The plot shows the effect of occlusion on the decision made by GoogleNet.

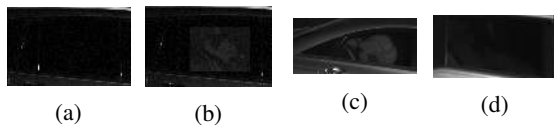


Figure 6: Effect of Tint and Illumination: (a) Tinted window (b) Histogram equalized. Child-safety seats: (c) Occluded (d) Non-occluded

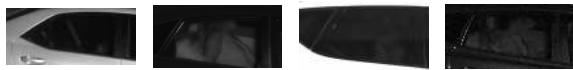


Figure 7: (a) No Person (b) Two Persons (c) Two Persons (d) Two Persons. Corner cases correctly predicted by VPDS. Faces have been redacted to preserve privacy.

Application Development, Deployment and Maintenance

We use a database of around 35000 raw images of front and rear views of vehicles. This dataset is then split randomly into 3 sets - train, validation and test which contain roughly 24000, 4000 and 7000 images respectively. The train set images are shown to the network while the validation set is used to do the model selection. We evaluate the performance of the model on the test set. We trained our models on an NVIDIA Tesla K80 GPU. We used DarkNet and Caffe for training YOLOv3 and the classifier models respectively.

VPDS has been deployed at more than three sites and functions efficiently at a very high accuracy. The system can be tuned based on the site requirements for more sensitive violator detection or less sensitive non-violator ticketing. The system has been successful in ticketing the 80% of violators arising from the voluntary compliance system and proved very beneficial for the transportation law enforcement agencies. VPDS requires a one time installation of the imaging equipment at the site. The images from the site are collected and a site-specific model is trained to achieve the best possible accuracy. Once the training is done there is minimal or no requirement of any intervention at the site except in

the rare event of any hardware failure. The local AI processing station is also configured one time and does not require monitoring except during a system upgrade. Also, the images collected during the process can be backed up elsewhere to avoid flooding the processing server.

Conclusion and Future Work

The reliability of voluntary compliance is questionable as studies have shown that 80% of the vehicles in an unmonitored HOV lane are in violation of the law. With the ever spreading urban sprawl and an overwhelming dependency of US cities on automobiles, decongestion is one of the highest priorities. We have developed VPDS, an AI based vehicle passenger detection system to effectively enforce HOV/HOT lane movement. VPDS automates and improves identification of HOV violators and assigns fines and tolls to HOV lane users. Moreover, it is extremely fast and takes less than 2s for classifying a vehicle as a violator or a non-violator with 96% accuracy without thwarting the normal traffic flow and using minimal hardware. Over a period of 2 years during which VPDS was deployed at three different sites, it has served approximately 30 million passengers. Serving roughly 1800 vehicles in morning and 2300 vehicles in the evening rush hours at one particular site, VPDS achieved an accuracy between 94-96% irrespective of the traffic flow or time of the day. This exemplifies an AI-based system which is highly accurate, consistent, fast, responsive in real-time, robust to externalities and requires little maintenance. In future, we aim to further improve the efficiency of VPDS by exploring recent advancements in object detection methods for better ROI detection and more powerful neural architectures for better classification. We also envision to make a holistic system with vehicle type identification as a sub-module working in conjugation with VPDS to automatically generate toll/fine so that congestion and violation management happen in a seamless manner.

References

Artan, Y., and Paul, P. 2013. Occupancy detection in vehicles using fisher vector image representation. *arXiv preprint arXiv:1312.6024*.

- Artan, Y.; Bulan, O.; Loce, R. P.; and Paul, P. 2016. Passenger compartment violation detection in hov/hot lanes. *IEEE Trans. Intelligent Transportation Systems* 17(2):395–405.
- Bulan, O.; Kozitsky, V.; Ramesh, P.; and Shreve, M. 2017. Segmentation-and annotation-free license plate recognition with deep localization and failure identification. *IEEE Transactions on Intelligent Transportation Systems* 18(9):2351–2363.
- Daley, W.; Arif, O.; Stewart, J.; Wood, J.; Usher, C.; Hanson, E.; Turgeson, J.; Britton, D.; et al. 2011. Sensing system development for hov/hot (high occupancy vehicle) lane monitoring.
- Fan, Z.; Islam, A. S.; Paul, P.; Xu, B.; and Mestha, L. K. 2013. Front seat vehicle occupancy detection via seat pattern recognition. US Patent 8,611,608.
- Felzenszwalb, P. F.; Girshick, R. B.; McAllester, D.; and Ramanan, D. 2010. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32(9):1627–1645.
- Hao, X.; Chen, H.; and Li, J. 2006. An automatic vehicle occupant counting algorithm based on face detection. In *Signal Processing, 2006 8th International Conference on*, volume 3. IEEE.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
2017. INRIX Global Traffic Scorecard. <http://inrix.com/scorecard/>. [Online; accessed 4-Sep-2018].
- Perronnin, F., and Dance, C. 2007. Fisher kernels on visual vocabularies for image categorization. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 1–8. IEEE.
- Perronnin, F.; Sánchez, J.; and Mensink, T. 2010. Improving the fisher kernel for large-scale image classification. In *European Conference on Computer Vision*, 143–156. Springer.
- Redmon, J., and Farhadi, A. 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.
- Schijns, S., and Mathews, P. 2005. A breakthrough in automated vehicle occupancy monitoring systems for hov/hot facilities. In *12th HOV Systems Conference*, volume 1.
- Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Sindagi, V. A., and Patel, V. M. 2018. A survey of recent advances in cnn-based single image crowd counting and density estimation. *Pattern Recognition Letters* 107:3–16.
2018. Number of cars in the US. <https://www.statista.com/statistics/183505/number-of-vehicles-in-the-united-states-since-1990/>. [Online; accessed 4-Sep-2018].
- Szegedy, C.; Liu, W.; Jia, Y.; Sermanet, P.; Reed, S.; Anguelov, D.; Erhan, D.; Vanhoucke, V.; and Rabinovich, A. 2015. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1–9.
- Wang, Y. R.; Xu, B.; and Paul, P. 2015. Determining a pixel classification threshold for vehicle occupancy detection. US Patent 9,202,118.
- Wshah, S.; Xu, B.; Bulan, O.; Kumar, J.; and Paul, P. 2016. Deep learning architectures for domain adaptation in hov/hot lane enforcement. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, 1–7. IEEE.
- Xu, B.; Paul, P.; Artan, Y.; and Perronnin, F. 2014. A machine learning approach to vehicle occupancy detection. In *Intelligent Transportation Systems (ITSC), 2014 IEEE 17th International Conference on*, 1232–1237. IEEE.
- Xu, B.; Paul, P.; and Perronnin, F. 2017. Vehicle occupancy detection using passenger to driver feature distance. US Patent 9,760,783.