



Cover Page



CLASSIFICATION OF OBJECTS USING CNN-BASED VISION AND LIDAR FUSION IN AUTONOMOUS VEHICLE ENVIRONMENT

¹G. Komali and ²Dr A Sri Nagesh

¹M. Tech Scholar and ²Professor

^{1&2}C.S. E, RVR&JC College of Engineering, Acharya Nagarjuna University,
Guntur, Andhra Pradesh, India

Abstract

In the past decade, Autonomous Vehicle Systems (AVS) have advanced at an exponential rate, particularly due to improvements in artificial intelligence, which have had a significant impact on social as well as road safety and the future of transportation systems. The fusion of light detection and ranging (LiDAR) and camera data in real-time is known to be a crucial process in many applications, such as in autonomous driving, industrial automation and robotics. Especially in the case of autonomous vehicles, the efficient fusion of data from these two types of sensors is important to enabling the depth of objects as well as the classification of objects at short and long distances. This paper presents classification of objects using CNN based vision and Light Detection and Ranging (LIDAR) fusion in autonomous vehicles in the environment. This method is based on convolutional neural network (CNN) and image up sampling theory. By creating a point cloud of LIDAR data up sampling and converting into pixel-level depth information, depth information is connected with Red Green Blue data and fed into a deep CNN. The proposed method can obtain informative feature representation for object classification in autonomous vehicle environment using the integrated vision and LIDAR data. This method is adopted to guarantee both object classification accuracy and minimal loss. Experimental results show the effectiveness and efficiency of presented approach for objects classification.

Keywords: Autonomous Vehicle Systems (AVS), Light Detection and Ranging (LIDAR), CNN.

Introduction

As technology constantly evolves, autonomous vehicles are becoming more popular, accessible, and affordable for more people in different countries and from different economic classes. Increasing accessibility results in a safer transportation experience, fewer deaths, and minimal injuries due to human-made mistakes that cause catastrophic accidents. To ensure the safety of individuals, it is necessary to deploy highly efficient and accurate learning models trained on a broad range of driving scenarios to precisely detect the surrounding objects under different weather and lighting conditions. This learning procedure via training will adjust the vehicle's decision-making process and control mechanism to take the necessary actions [1].

The interest in autonomous vehicles has increased in recent years due to the advances in multiple engineering fields such as machine learning, robotic systems and sensor fusion. The progress of these techniques leads to more robust and trustworthy computer vision algorithms. Using sensors such as Laser Imaging Detection and Ranging (LiDAR), radar, camera or ultrasonic sensors with these techniques enables the system to detect relevant targets in highly dynamic surrounding scenarios. These targets may include pedestrians, cyclists, cars or motorbikes among others, as discussed in public autonomous car datasets [2].

Sensors in autonomous vehicles (AV) are used in two main categories: localization to measure where the vehicle is on the road and environment perception to detect what is around the vehicle. Global Navigation Satellite System (GNSS), Inertial Measurement Unit (IMU), and vehicle odometry sensors are used to localize the AV. Localization is needed, so the AV knows its position with respect to the environment. Autonomous vehicles should be instantaneous, accurate, stable, and efficient in computations to produce safe and acceptable traveling trajectories in numerous urban to suburb scenarios and from high-density traffic flow to high-speed highways. In real-world traffic, various uncertainties and complexities surround road and weather conditions, whereas a dynamic interaction exists between objects and obstacles, and tires and driving terrains. An autonomous vehicle must rapidly and accurately detect, recognize, and classify and track dynamic objects with complex backgrounds and posing technical challenges.

Driving environment recognition serves as a driving environment dynamic, static object detection, lane detection, and vehicle location estimation based on sensors that can obtain information about the driving environment, and the decision to determine the vehicle trajectory, such as the creation and avoidance of routes to the destination. Longitudinal and lateral controls are performed to reliably drive the target control values of the vehicle determined by recognition and decision [3].



Cover Page



To ensure correct and safe driving, a fundamental pillar in the Autonomous Driving (AD) system is perception, which leverages sensors such as cameras and LiDARs (Light Detection and Ranging) to detect surrounding obstacles in real time [4].

Autonomous vehicles rely on their perception systems to acquire information about their immediate surroundings. It is necessary to detect the presence of other vehicles, pedestrians, and other relevant entities. Safety concerns and the need for accurate estimations have led to the introduction of lidar systems to complement camera- or radar-based perception systems [5].

Regarding object distance estimation, there are several approaches that have been proposed, depending on the modality of the sensors used, such as radar, LiDAR, or camera. Each sensor modality is capable of perceiving the environment with a specific perspective and is limited by detecting certain attribute information of objects. More specifically, vision-based approaches are more robust and accurate in object detection but fail in estimating the distance of the object accurately [6].

Deep learning algorithms have been utilized in different aspects of AV systems, such as perception, mapping, and decision making. These algorithms have proven their ability to solve many of these difficulties, including computational loads faced by traditional algorithms while maintaining decent accuracy and fast processing speed. Currently high-performance vision system usually is based on deep learning techniques. Deep neural networks (DNN) have proven to be an extremely powerful tool for many vision tasks. Hence in this paper classification of objects using CNN based Vision and LIDAR fusion in autonomous vehicle environment is presented.

literature survey:

G Ajay Kumar, Jin Hee Lee, Jongrak Hwang, Jaehyeong Park, Sung Hoon Youn and Soon Kwon et. al., [7] presents LiDAR and Camera Fusion Approach for Object Distance Estimation in Self-Driving Vehicles. This paper presents a method to estimate the distance (depth) between a self-driving car and other vehicles, objects, and signboards on its path using the accurate fusion approach. Based on the geometrical transformation and projection, low-level sensor fusion was performed between a camera and LiDAR using a 3D marker. Jian Nie, Jun Yan, Huilin Yin, Lei Ren, and Qian Meng et. Al [8] presents A Multimodality Fusion Deep Neural Network and Safety Test Strategy for Intelligent Vehicles. In this paper, they firstly propose a multimodality fusion framework called Integrated Multimodality Fusion Deep Neural Network (IMF-DNN), which can flexibly accomplish both object detection and end-to-end driving policy for prediction of steering angle and speed.

Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Won Park, Sara Rampazzi, Qi Alfred Chen, Z. Morley Mao, Kevin Fu et. al., [9] presents Adversarial Sensor Attack on LiDAR-based Perception in Autonomous Driving. we perform the first security study of LiDAR-based perception in AV settings, which is highly important but unexplored. We consider LiDAR spoofing attacks as the threat model and set the attack goal as spoofing obstacles close to the front of a victim AV.

Mhafuzul Islam, Mashrur Chowdhury, Hongda Li, and Hongxin Hu et. al., [10] presents Vision-Based Navigation of Autonomous Vehicles in Roadway Environments with Unexpected Hazards. They develop a DNN-based autonomous vehicle driving system using object detection and semantic segmentation to mitigate the adverse effect of this type of hazard, which helps the autonomous vehicle to navigate safely around such hazards. We find that our developed DNN-based autonomous vehicle driving system, including hazardous object detection and semantic segmentation.

Babak Shahian Jahromi, Theja Tulabandhula and Sabri Cetin et. al., [11] presents Real-Time Hybrid Multi-Sensor Fusion Framework for Perception in Autonomous Vehicles. We propose a new hybrid multi-sensor fusion pipeline configuration that performs environment perception for autonomous vehicles such as road segmentation, obstacle detection, and tracking Tested on over 3K road scenes, our fusion algorithm shows better performance in various environment scenarios compared to baseline benchmark networks.

Bike Chen, Chen Gong, Jian Yang et. al., [12] discussed about Importance-Aware Semantic Segmentation for Autonomous Vehicles. The IAL (Importance Aware Losses) operates under a hierarchical structure and the classes with different importance are located in different levels so that they are assigned distinct weights They derive the forward and backward propagation rules for IAL and apply them to four typical deep neural networks for realizing SS in an intelligent driving system.

Jin Fang, Feilong Yan, Tongtong Zhao, Feihu Zhang, Dingfu Zhou, Ruigang Yang, Yu Ma and Liang Wang et. al., [13] presents Simulating LIDAR Point Cloud for Autonomous Driving using Real-world Scenes and Traffic Flows. We present a LIDAR simulation framework that can automatically generate 3D point cloud based on LIDAR type and placement.

Xinxin Du, Marcelo H. Ang Jr. and Daniela Rus et. al., [14] presents Car Detection for Autonomous Vehicle: LIDAR and Vision Fusion Approach Through Deep Learning Framework. They propose a LIDAR and vision fusion system for car detection through the deep learning framework. With further optimization of the framework structure, it has great potentials to be implemented onto the autonomous vehicle.

Andreas Eitel Jost Tobias Springenberg Luciano Spinello Martin Riedmiller Wolfram Burgard et. al., [15] presents Multimodal Deep Learning for Robust RGB-D Object Recognition. architecture is composed of two separate CNN processing streams: one for each modality which are consecutively combined with a late fusion network. They focus on learning with imperfect sensor data, a typical problem in real-world robotics tasks.

CLASSIFICATION OF OBJECTS USING CNN.

In this work, classification of objects using CNN based vision and LIDAR fusion in autonomous vehicle environment is presented. The framework of presented model is shown in Fig. 1.

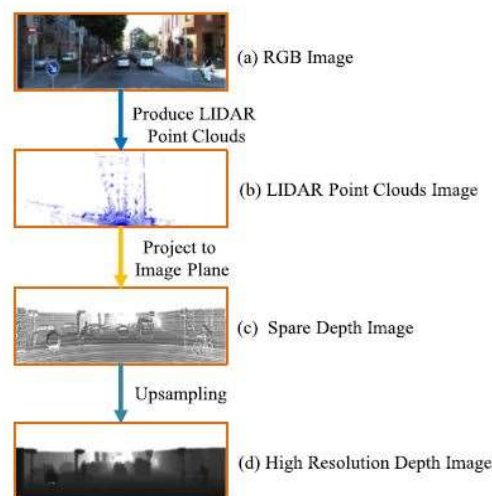


Fig. 1: THE FRAMEWORK OF PRESENTED MODEL

Autonomous vehicles use various sensors, such as LiDAR, radar, cameras, and ultrasonic sensors, to map and recognize the environment surrounding the vehicle. While considering real-time performance, object detection, classification ability, and accurate distance estimation, LiDAR and vision fusion techniques have been introduced for object detection based on different levels of data fusion.

A Laser Imaging Detection and Ranging (LiDAR) sensor has been used in this project to gather information regarding the environment. Differently from the camera, the LiDAR sensor transmits laser pulses and measures the time it takes until a reflection is received. Based on this, it calculates the distance of the target and its 3D coordinates since the angles used to send the laser pulse and the distance are known.

One of the largest and widely used benchmark datasets in the autonomous driving research community is the KITTI dataset and is used here which provides LiDAR point clouds, stereo color, and grayscale pictures, and GPS coordinates. The data was captured on the highways and rural areas of a mid-sized city in Germany called Karlsruhe. The tasks that can utilize this dataset include 3D Object Detection, Visual Odometry, Stereo Matching, and Optical Flow. The Object Detection part of the dataset consists of 7,481 training and 7,518 test images, with annotated boxes around the objects of interest.



Cover Page



The camera and LIDAR on the vehicle are used to collect images. We first capture the sparse-depth map by rotating Velodyne laser-point cloud data from the KITTI database to the RGB image plane using the calibration matrix. Then, we upsample the sparse depth map to high-resolution depth image. We extract four objects (pedestrian, cyclist, car, and truck) from each image by considering the ground truth from KITTI. The rapid growth of research and commercial enterprises relating to autonomous robots, drones, humanoid robots, and AVs has established a high demand for LiDAR sensors due to its performance attributes such as measurement range and accuracy, robustness to surrounding changes and high scanning speed.

We build three image datasets according to these objects. One database is for the pure RGB image of the four kinds of object, one for the gray-scale image with gray level corresponding to actual distance information from LIDAR point clouds, and the third one is an RGB-LIDAR image dataset consisting of the former two information. Each dataset comprises 6843 labeled objects. Finally, we present a structure based on CNN to train a classifier for detecting the four kinds of objects on the road. These classification results are provided to the driving cognitive module for vehicle decision-making and control.

LiDAR gives us rich information with point clouds (which include position coordinates x, y, z and intensity i) as well as a depth map. The first step is to process the LiDAR point cloud data (PCD). The raw LiDAR point cloud is high resolution and covers a long distance (for example a 64-lens laser acquires more than 1 million points per second). Dealing with such a large number of points is very computationally expensive and will affect the real-time performance of our pipeline.

In this study, a novel method of up sampling LIDAR range inputs is employed to align depth with RGB images. In this method, we compute dense depth maps just from the original range data instead of using information from RGB images. We formulate up-sampling using bilateral filtering formalism in our method to generate the dense map D (output image) from a noisy and sparse-depth image I. Assuming that input I is coordinated in pixel units and features calibration w.r.t. a high-resolution camera, pixel positions in I are nonintegers owing to the uncertainty of calibration parameters and data sparsity. According to the intensity value of a pixel p on the depth map, expressed as lower index (P and its N neighborhood mask, the pixel value lies on the same position of output map Dp , as shown in the following equation:

$$D_p = \frac{1}{W_p} \sum_{q \in N} G_{\sigma_r}(|I_q|) I_q G_{\sigma_s}(\|p - q\|) \quad (1)$$

Where G_{σ_r} penalizes the influence of points q caused by their range values, G_{σ_s} weighs inversely to the distance between position p and location q, and W_p functions as the normalized factor, which ensures that the sum of weights are equal to one. In (1), we set G_{σ_s} to be inversely proportional to the Euclidean distance ($\|p - q\|$) between pixel position p and location q.

RGB images and the 3-D point clouds from KITTI are used as object benchmarks to classify objects, such as cars, pedestrians, trucks, and cyclists. RGB color images are captured by the left color video camera (10 Hz, resolution: 1392 × 512 pixels, opening: 90° × 35°), whereas the 3-D point clouds are produced by a Velodyne HDL-64E unit and projected back in image forms. As one of the few available sensors that provide depth information, Velodyne system can generate accurate 3-D data from moving platforms. This system can also be applied in outdoor scenarios and long sensing range compared with structured light systems such as Microsoft Kinect.

Convolutional neural networks (CNN) have been extensively applied to image classification and computer vision, and have returned 100% classification rates on datasets such as ImageNet. In CNN architecture, successive layers of neurons learn progressively complex features in a supervised way by back-propagating classification errors, with the last layer representing output image categories. CNNs do not use a distinct feature extraction module or a classification module, i.e. CNNs do not have an unsupervised pre-training and the input representation is implicitly through supervised training but eliminate the need for manual feature description and feature extraction. It extracts features from raw data based on pixel values leading to final object categories. Each layer in the CNN finds successively complex features where the first layer finds a small, simple feature anywhere on the image, the second layer finds more complex features and so on. At the last layer, these feature maps are processed using fully connected neural networks (FCNN).

AlexNet is a leading architecture for any object-detection task and may have huge applications in the computer vision sector of artificial intelligence problems. In the future, AlexNet may be adopted more than CNNs for image tasks. The Alexnet has eight

layers with learnable parameters. The model consists of five layers with a combination of max pooling followed by 3 fully connected layers and they use ReLu (Rectified Linear Units) activation in each of these layers except the output layer.

For object classification, we classified images from KITTI into cars, cyclists, pedestrians, and trucks. Then, we adopt the AlexNet model as our CNN architecture. AlexNet comprises five convolutional layers (named conv1–conv5) and three fully connected layers (named as fc6, fc7, and 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 comprises multiple neurons, and each neuron contains a positive value and is connected to all neurons in the previous layer. We resize the captured images to 128×128 resolution for valid input and then passed them into AlexNet. AlexNet is trained for 1000 classes. We change the size of fc8 layer from 1000 to 4 to match our dataset with four classes. The parameters from layer conv1 to layer fc6 are fixed to prevent overfitting.

This RGB-LIDAR-based method notably improves average precision on classification of four categories in the KITTI datasets. We use the same dataset for training and testing models.

RESULT ANALYSIS

In this section the design and implementation results of classification of objects using CNN based fusion of vision and LIDAR in autonomous vehicle environment are discussed.

To implement this project we have designed following modules: i) Upload Kitti Dataset: using this module we will upload dataset to application, ii)Load Alexnet LIDAR CNN Model: This module will read all images and then applying up-sampling to increase image intensity and then extract RGB value to train ALEXNET CNN model, iii) Run LIDAR Object Detection & Classification: In this module we we will upload test image and then Alexnet model will detect and classify objects from that test image and iv) LIDAR Accuracy & Loss Graph: Using this module we will plot LIDAR Alexnet accuracy and loss graph and we train this algorithm for 10 EPOCH and in graph we will get accuracy plot for each epoch.

Click on ‘Upload Kitti Dataset’ button to upload dataset to application.

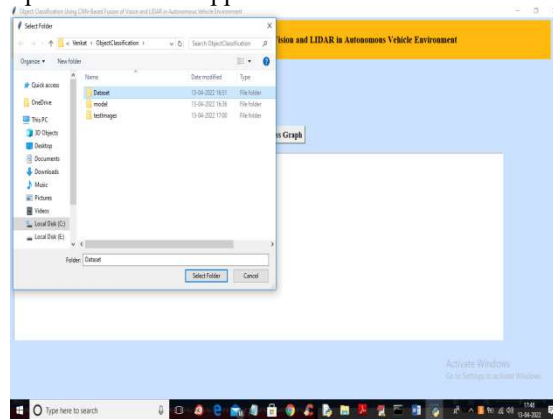


Fig. 2: KITTI DATASET SCREEN

In above screen selecting and uploading dataset folder and this folder contains different types of objects. The dataset has 3 different types of classes and just go inside any folder to view those type of images. The loaded dataset is shown in Fig. 3.

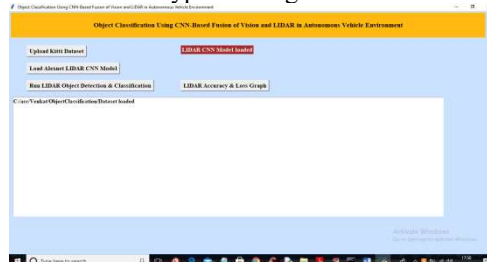


Fig. 3: LOADED DATASET SCREEN

In above screen in red colour text we can see CNN model loaded and now click on ‘Run LIDAR Object Detection & Classification’ button to upload test image and get below output.

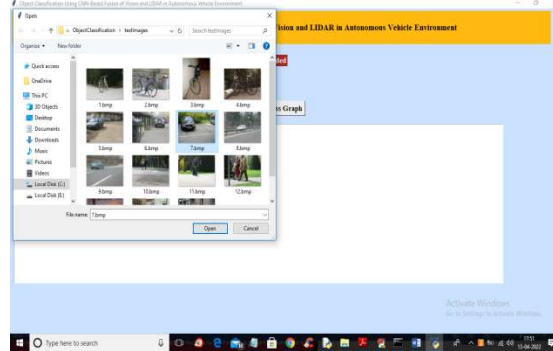


Fig. 4: TEST IMAGE UPLOADING SCREEN

In above screen, selecting and uploading 7.bmp im and then click on ‘Open’ button to get below output.

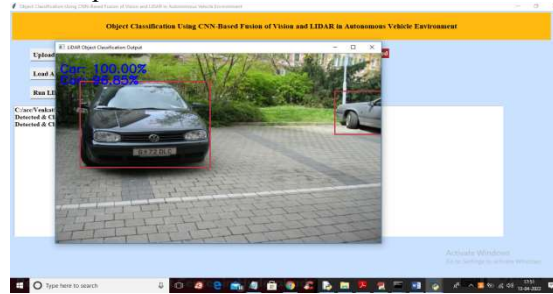
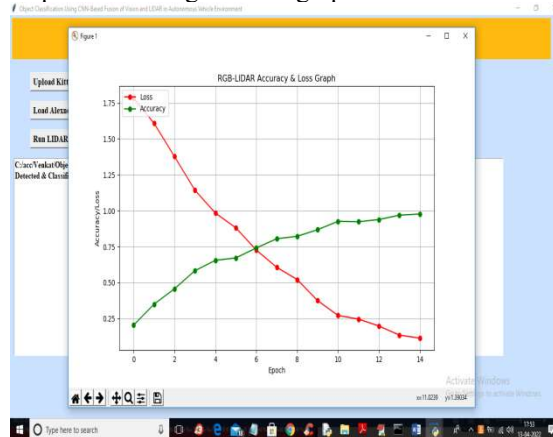


Fig. 5: CLASSIFIED OBJECT IMAGE

now click on ‘LIDAR Accuracy & Loss Graph’ button to get below graph.



6: FINAL OUTPUT GRAPH

CONCLUSION

In this work, classification of objects using CNN based vision and LIDAR fusion in autonomous vehicle environment is presented. we propose a deep-learning-based approach by fusing vision and LIDAR data for object detection and classification in autonomous vehicle environment. On the one hand, we upsample point clouds of LIDAR data and convert the upsampled point cloud data into pixel-level depth featuremap. On the other hand, we convert the RGB together with depth feature map and then fed the data into a CNN. On the basis of the integrated RGB and depth data, we utilize DCNN to perform feature learning from raw input



Cover Page



information and obtain informative feature representation to classify objects in the autonomous vehicle environment. The proposed approach, in which visual data are fused with LIDAR data, exhibits superior classification accuracy over the approach using only RGB data or depth data. During the training phase, using LIDAR information can accelerate feature learning and hasten the convergence of CNN on the target task. We perform experiments using the public dataset and display the effectiveness and efficiency of the proposed approach.

References

1. Hrag-Harout Jebamikyous, (Member, Ieee), And Rasha Kashef, “Autonomous Vehicles Perception (AVP) Using Deep Learning: Modeling, Assessment, and Challenges”, IEEE ACCESS, VOLUME 10, 2022, doi: 10.1109/ACCESS.2022.3144407
2. Javier Mendez, Miguel Molina, Noel Rodriguez, Manuel P. Cuellar and Diego P. Morales, “Camera-LiDAR Multi-Level Sensor Fusion for Target Detection at the Network Edge”, Sensors 2021, 21, 3992, doi.org/10.3390/s21123992
3. Mingyu Park, Hyeonseok Kim and Seongkeun Park, “A Convolutional Neural Network-Based End-to-End Self-Driving Using LiDAR and Camera Fusion: Analysis Perspectives in a Real-World Environment”, Electronics 2021, 10, 2608, doi.org/10.3390/electronics10212608
4. Yulong Cao, Ningfei Wang, Chaowei Xiao, Dawei Yang, Jin Fangz Ruigang Yanggy Qi Alfred Chen, Mingyan Liux Bo Li, “Invisible for both Camera and LiDAR: Security of Multi-Sensor Fusion based Perception in Autonomous Driving Under Physical-World Attacks”, arXiv:2106.09249v1 [cs.CR] 17 Jun 2021
5. G Ajay Kumar, Jin Hee Lee, Jongrak Hwang, Jaehyeong Park, Sung Hoon Youn and Soon Kwon, “LiDAR and Camera Fusion Approach for Object Distance Estimation in Self-Driving Vehicles”, Symmetry 2020, 12, 324; doi:10.3390/sym12020324
6. You Li and Javier Ibanez-Guzman, “Lidar for Autonomous Driving”, IEEE SIGNAL PROCESSING MAGAZINE, 1053-5888/20©2020IEEE, Digital Object Identifier 10.1109/MSP.2020.2973615
7. G Ajay Kumar, Jin Hee Lee, Jongrak Hwang, Jaehyeong Park, Sung Hoon Youn and Soon Kwon, “LiDAR and Camera Fusion Approach for Object Distance Estimation in Self-Driving Vehicles”, Symmetry 2020, 12, 324; doi:10.3390/sym12020324
8. Jian Nie, Jun Yan, Huilin Yin, Lei Ren, and Qian Meng, “A Multimodality Fusion Deep Neural Network and Safety Test Strategy for Intelligent Vehicles”, 2020 IEEE Transactions on Intelligent Vehicles, 2379-8858 (c), DOI 10.1109/TIV.2020.3027319
9. Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Won Park, Sara Rampazzi, Qi Alfred Chen, Z. Morley Mao, Kevin Fu, “Adversarial Sensor Attack on LiDAR-based Perception in
10. Autonomous Driving”, CCS '19, November 11–15, 2019, London, United Kingdom, 2019 Association for Computing Machinery, doi.org/10.1145/3319535.3339815
11. Mhafuzul Islam, Mashrur Chowdhury, Hongda Li, and Hongxin Hu, “Vision-Based Navigation of Autonomous Vehicles in Roadway Environments with Unexpected Hazards”, Transportation Research Board 2019, DOI: 10.1177/0361198119855606