

Cooper: Cooperative Perception for Connected Autonomous Vehicles based on 3D Point Clouds

Qi Chen*, Sihai Tang*, Qing Yang[†] and Song Fu[†]

*{QiChen, SihaiTang}@my.unt.edu, [†]{Qing.Yang, Song.Fu}@unt.edu

Department of Computer Science and Engineering
University of North Texas
Denton, TX, USA

Abstract—Autonomous vehicles may make wrong decisions due to inaccurate detection and recognition. Therefore, an intelligent vehicle can combine its own data with that of other vehicles to enhance perceptive ability, and thus improve detection accuracy and driving safety. However, multi-vehicle cooperation perception requires physical integration of real world scenes and the traffic of data exchange far exceeds the bandwidth of existing vehicular networks such as Dedicated short-range communication (DSRC). To the best of our knowledge, we are the first to conduct a study on cooperation perception towards the effort of enhancing the detection ability of self-driving systems. In this work, relying on 3D point clouds, sensor data from different positions and angles of connected vehicles are merged together. Specifically, the merging and perception of point cloud with different resolution/density is taken into account. A point cloud based 3D objects detection is proposed to work on a diversity of aligned point clouds. Evaluated on KITTI and our proposed T&J dataset, it shows that the proposed system outperforms individual perception in the fields of extending sensing area, improving detection accuracy and promoting augmented results. We find that collaboration offers comprehensive information as opposed to information perceived by individuals alone. Most importantly, we demonstrate that the bandwidth of DSRC can satisfy point clouds transmission for cooperative perception.

I. INTRODUCTION

A significant part of the push towards autonomous driving vehicles, or self-driving vehicles, has been supported by the prospect that they will save lives by getting involved in fewer crashes with fewer injuries and deaths than human-driven cars. However, up until this point, most comparisons between human driven cars and self-driving vehicles have been unbalanced and contain various unfair elements. Self-driving cars do not experience fatigue, emotional debilitations such as anger or frustration. But, they are unable to react to uncertain and ambiguous situations with the same skill or anticipation of an attentive and seasoned human driver.

Similarly, isolated self driving vehicles may make wrong decision due to the failure of objects detection and recognition. Just as a human driver will make bad decisions while under the influence, such decisions made by the vehicle based on these failures will prove just as bad or worse than their human counterpart. Such vehicles must completely rely on itself for decision making, and thus will not have the privilege of data redundancy, i.e., no information is received from nearby vehicles. Sensor failure or any other technical error will lead to fallacious results, leading to disastrous impacts.

A. Motivations

The deficit of data due to single source will ultimately have a negative impact as well. Take the example of Tesla's crash in California, the car made a fatal decision because its sensors picked up the concrete barrier but discarded the information due its immobile state on the radar. Similarly, the Tesla crash against a fire truck in Los Angeles had the same issue of making a bad decision on limited sensor data [25]. Of course, there are also instance of various other circumstances leading up to bad decisions, such as the Uber training incident [16]. In this case, the vehicle did detect an unknown object, the pedestrian, from a distance. As the vehicle approached the unknown object, it gradually discerned the object to be a vehicle and finally a pedestrian, but by then, it was too late.

We further explore the reasons why detection failure happened. It is easy to determine that some detection failures are caused due to objects being blocked or existing in the blind zones of the sensors. Detection failures could also be caused by bad recognition because the received signal is too weak or because the signal is missing due to system malfunction.

Our motivation comes from these incidents, because in contrast to isolated autonomous driving vehicles, like the ones in the accidents, connected autonomous vehicles (CAV) can share their collected data with each other leading to more information. We propose that information sharing can improve driving performance and experiences. Constructive data redundancy will provide endless possibilities for safe driving and multiple vehicles can collaborate together to compensate for data scarcity and provide a whole new scope for the vehicle in need. Autonomous vehicles have powerful perception systems, and together, they can achieve a proper data sharing and analysis platform to gain much more reliability and accuracy[28].

B. Limitations of Prior Work

Although adding connectivity to vehicles has its benefits, it also has challenges. By adding connectivity, there can be issues with security, privacy, and data analytics and aggregation due to the large volume of information being accessed and shared.

Current state of multi-sensor fusion consists of three distinct categories: low level fusion, feature level fusion, and high level fusion [22]. Each of these categories possess its own unique advantages and disadvantages. As their names imply, low level fusion consists of raw data fusion without any pre-processing

done to the data. Feature-level fusion takes the features extracted from the raw data before fusion. Finally, high level fusion takes the objects detected from each individual sensors and conducts the fusion on the object detection results [22].

High level fusion is often opted over the other two levels of fusion due to being less complex, but this is not suitable for our needs. Object level relies too heavily on single vehicular sensors and will only work when both vehicles share a reference object in their detection. This does not solve the issue of previously undetected objects, which will remain undetected even after fusion. And thus, we turn our sights on the other two categories.

C. Proposed Solution

To tackle the issue, we look at one of the base categories, the low level fusion of raw data. Raw sensing data is an integral part of all sensors on autonomous driving vehicle, therefore, it is very suitable for transferring them between different cars from various manufactures. As such, the heterogeneity of different data processing algorithms would not affect the accuracy of the data being shared among vehicles. As autonomous driving is of and in itself a crucial task, being so integrated in the vehicle, even a single small error in detection can lead to a catastrophic accident. Therefore, we need the autonomous cars to perceive the environment with as much clarity as possible. To achieve this end goal, they will need a robust and reliable perception system.

Two major issues that we seek to address in doing so are as follows: (1) the type of data that we need to share among vehicles, and (2) the amount of the data that needs to be transferred versus the amount of data that is actually necessary to the recipient vehicle. The first issue arises with the shareable data within the dataset native to the car. The second problem exists in the sheer amount of data that each vehicle generates. Since each autonomous vehicle will collect more than 1000GB of data [2] every day the challenge of assembling only the regional data becomes even harder. Similarly, reconstructing the shared data collected from different positions and angles by nearby perception system is another major challenge.

Of the different types of raw data, we propose to use the LiDAR point clouds as a solution for the following reasons:

- LiDAR point clouds have the advantage of spatial dimension over 2D images and video.
- Native obfuscation of entities or private data such as people's faces and license plate numbers while preserving the accurate model of the perceived object.
- Versatility in the fusion process over images and video due to the data being consisted from points rather than pixels. For image or video fusion, the requirement is a clear zone of overlap, and this is unnecessary for point cloud data, making this a much more robust choice, especially when taking the different possible point of views of cars into perspective.

With the three different highlights of using the raw LiDAR data as our fusion substrate, we propose the **Cooperative**

Perception (Cooper) system for connected autonomous vehicles based on 3D point clouds.

D. Contributions

Inaccurate object detection and recognition are major impediments in achieving a powerful and effective perception system. Autonomous vehicle eventually succumb to this inability and fail to deliver the expected outcome, which is unsafe to autonomous driving. To address these issues we have proposed a solution in which an autonomous vehicle combines its own sensing data with that of other connected vehicles to help enhance perception. We also believe that data redundancy, as mentioned, is the solution to this problem and we can achieve it through data sharing and combination between autonomous vehicles. The proposed Cooper system can improve the detection performance and driving experience thus providing protection and safety.

To achieve our proposed solution, we will perform the following steps in our experiments:

- We propose a Sparse Point-cloud Object Detection (SPOD) methods using low-density point clouds. It can also work on high-density LiDAR data, which makes Cooper on multi-vehicles' possible.
- We show how the proposed system outperforms individual perception in the fields of extending sensing area and improving detection accuracy.
- we demonstrate that the bandwidth of DSRC can satisfy point clouds transmission for Cooper based on LiDAR point cloud data.

II. COOPERATIVE SENSING

Given the current outlook and work done in the field of fusion and information usage in autonomous vehicles, we need to go a step further and truly define what we see as cooperative sensing. We see this as a series of challenges and benefits that will be an unavoidable part of progress.

A. Benefits of Sharing

Based on our observations, we wonder if detection accuracy can be improved using sensor data from multiple cars. As we know, the sensing devices on autonomous vehicles work together to map the local environment and monitor the motion surrounding vehicles. According to the collected data, shareable resources can be extracted from these vehicles. For example, there is a blocked area region behind obstacles on the road that could not be sensed by one car but data gathered for this same area can be sensed and provided by other nearby cars. Meanwhile, vehicles on adjacent districts or crowded zones can keep connection for a longer duration, thereby enhancing cooperative sensing, which will greatly help other vehicles by providing crucial information. Hence, we propose a cooperative perception method to improve autonomous driving performance. This framework facilitates a vehicle to combine its sensor data with that of its cooperators' to enhance perceptive ability, and thus improving detection accuracy and driving safety.

B. Difficulty of Sharing

Even though shareable resources offer useful information, vehicles prefer to utilize raw data rather than extracted results. The detected results from other cars are hard to authenticate and trust issues further complicate this matter. Also, since sharing all collected data is also impractical, we need to take into consideration the bandwidth and latency of vehicular networks. First, the bandwidth and latency of vehicular networks must satisfy data transmission for cooperative perception. Then, the vehicles need to reconstruct the received data because it was taken on different positions and angles. With this series of questions, we elaborate our research on building cooperative perception.

C. Data Choice

First, we demonstrate which type of sensing data is suitable for cooperative perception. Noting that perception systems are mainly developed on image-based and LiDAR-based sensor data. As we mention before, image data holds advantage on object classification and recognition while lacking on location information. In the next section, our proposed SPOD method overcomes the shortcomings of point clouds, which were too sparse to detect objects. Based on the above reasons, we make a priority of these two sensor data for cooperative sensing. We prefer LiDAR data because it holds advantage in providing location information [21]. By only extracting positional coordinates and reflection value, point clouds can be compress into 200 KB per scan. For some applications, such as small object detection, for example license plate tracking, it is difficult for point clouds to recognize plate information. However, when utilized with cooperative perception, we are still able to locate the plates in point clouds and ask for its image data from connected vehicles. Because image and LiDAR point clouds are aligned together in perception system's installation, we integrate the above demand-driven strategy mainly relying on point clouds. In some cases, it is necessary to extract a fragment of the image data in cooperative perception.

D. Data reconstruction

Also, vehicles need to reconstruct the received data because it was taken on different positions and angles. By exchanging LiDAR data, local environment can be reconstructed intuitively by merging point clouds into its physical positions. In order to reconstruct local environment by mapping point clouds into physical positions, additional information is encapsulated into the exchange package. Said package should be constituted from LiDAR sensor installation information and its GPS reading, which determines the center point position of every frame of point clouds. Vehicle's IMU (inertial measurement unit) reading is also required because it records the offset information of the vehicle during driving: it represents a rotation whose yaw, pitch, and roll angles are α , β and γ , respectively [24]. A rotation matrix R will be generated in Equation 1.

$$R = R_z(\alpha)R_y(\beta)R_x(\gamma) \quad (1)$$

Here $R_z(\alpha)$, $R_y(\beta)$, $R_x(\gamma)$ are three basic rotation matrices rotate vectors by an angle on the z-, y-, x-axis in three dimensions:

$$R_z(\alpha) = \begin{bmatrix} \cos\alpha & -\sin\alpha & 0 \\ \sin\alpha & \cos\alpha & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R_y(\beta) = \begin{bmatrix} \cos\beta & 0 & \sin\beta \\ 0 & 1 & 0 \\ -\sin\beta & 0 & \cos\beta \end{bmatrix}$$

$$R_x(\gamma) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\gamma & -\sin\gamma \\ 0 & \sin\gamma & \cos\gamma \end{bmatrix}$$

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} X_R \\ Y_R \\ Z_R \end{bmatrix} \cup \begin{bmatrix} X'_T \\ Y'_T \\ Z'_T \end{bmatrix} \quad (2)$$

$$\begin{bmatrix} X'_T \\ Y'_T \\ Z'_T \end{bmatrix} = R \times \begin{bmatrix} X_T \\ Y_T \\ Z_T \end{bmatrix} + \begin{bmatrix} \Delta d_{x_T} \\ \Delta d_{y_T} \\ \Delta d_{z_T} \end{bmatrix} \quad (3)$$

When connected vehicles exchange message, cooperative perception produces a new frame by combining transmitter and receiver's sensor data using Equation 2, where we have the set of all coordinates equal to the coordinates of the receiver union with the the coordinates from the transmitter. However, as the transmitting vehicle is in a different state than the receiver, we must apply a transform to the original coordinates so that they match the state of the receiving vehicle.

To obtain the correct state for the transmitter's orientation, we use Equation 1.

Note $[X'_T Y'_T Z'_T]'$ is the transmitter's point cloud after applying the transform R to the translated coordinates of the transmitting vehicle. The transform is calculated by Equation 1, using the IMU value difference between the transmitter and the receiver.

III. COOPERATIVE PERCEPTION

To detect objects from the fused point cloud, we must consider a few factors. Most LiDAR being used generate three known types of density clouds, produced by Velodyne [3] LiDAR devices, being 64-beam, 32-beam and 16-beam. However, as we cannot state that all vehicles will use the top model of 64-beam, we must consider how sparse data point cloud will impact our method so that our method will work for a broad spectrum of density clouds.

A. Object Detection based on Point Clouds

As we know, each self-driving vehicle will extract sensor data to perceive details in the local environment, such as lane detection, traffic sign detection and objects like cars, cyclists and pedestrians. However, accurate detection of objects in point clouds is a challenge due to LiDAR point clouds being sparse and it having a highly variable point density. For example, recently, based on point clouds dataset in KITTI [8], VoxelNet [29] has announced its experiments on car detection task which outperformed the state-of-the-art 3D detection

methods. Its car detection average precision is 89.60%, and for smaller objects, such as pedestrians and cyclists, the average precision drops to 65.95% and 74.41% respectively in a fully visible (easy) detecting environment. While in a difficult to see (hard) detecting condition, the car, pedestrian and cyclist detection further drop to 78.57%, 56.98%, and 50.49%, respectively. Another insight here is that LiDAR provides sparse 3D point clouds with location information but is hard to classify and recognize. To analyze the results from the above works, we cannot ignore the failure detection. This allows us to approach the issue from another perspective - cooperative sensing methods to improve detection accuracy.

B. Sparse Point-cloud Object Detection (SPOD)

Typically autonomous vehicles use single end-to-end deep neural network to operate on a raw point cloud. However, after cooperative sensing, the re-constructed data from different LiDAR devices may have different features like point density. For example, Velodyne [3] produces 64-beam, 32-beam and 16-beam LiDAR devices, which provide different density point clouds. Similar to image's resolution, 3D detector using deep neural network may have inaccuracy recognition results when used on low density point clouds. We note that 64-beam LiDAR, which provide the highest resolution LiDAR data, is well adopted by researches and companies on 3D object detection [29], [27]. While some others, as in our case, use 16-beam LiDAR, which outputs sparse data but has a price advantage over its higher end counterparts. This requires our proposed detection method on its assembled 3D detection model not only to work on high density data, but also can detect objects from much sparser point clouds. Unfortunately, these convolutional neural network (CNN)-based object detection methods are not suitable for low-density data because of insufficient of input features. Inspired by the state-of-the-art work [27], we propose the Sparse Point-cloud Object Detection (SPOD) methods which can adapt low density point clouds.

C. Architecture of SPOD

The proposed detector, depicted in Fig. 1, consists of three components. Our adopted 3D LiDAR point cloud is represented as a set of cartesian coordinates, (x, y, z) with reflection values. The distribution of point clouds is much too sparse and irregular. Specifically in the preprocessing, to obtain a more compact representation, point clouds are projected onto a sphere using approach from [26] to generate a dense representation. In voxel feature extractor components, our framework takes represented point clouds as input, feeding extract voxel-wise features to voxel feature encoding layer, this is well demonstrated by Voxelnet [29]. Then a sparse convolutional middle layer [14] is applied. Sparse CNN offers computational benefits in LiDAR-based detection because the grouping step for point clouds will generate a large number of sparse voxels. In this approach, output points are not computed if there is no related input points. Finally, Region Proposal Network (RPN) [20] is constructed using single shot multibox

detector (SSD) architecture [15]. The feature maps as input to RPN from Sparse CNN and are concatenate into one feature map for prediction. Framework in every vehicle use this single end-to-end trainable network to produce 3D detection results not only from dense LiDAR data but also from low resolution LiDAR data from nearby vehicles.

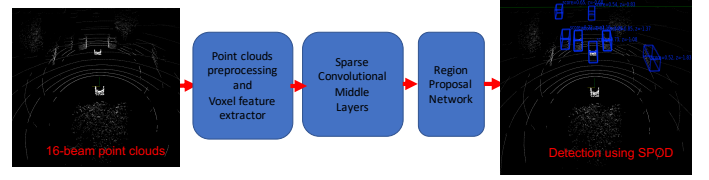


Fig. 1: 3D object detector works on 16-beam raw point clouds

Eventually, we successfully adopt SPOD to detect objects both on our collected sparse data and on dense KITTI data. In the next section, we demonstrate a full evaluation of SPOD detection.

IV. EVALUATION AND RESULT ANALYSIS

In this section, we evaluate the performance of our proposed cooperative perception using real world sensor data.

A. Datasets

In the experiment, we test Cooper based on two datasets, KITTI and T&J, thereby obtaining two types of density point clouds, dense vs sparse. In the dense KITTI dataset, a 64-beam LiDAR sensor is used to collect point clouds. But in our T&J dataset, which supplies 16-beam point cloud, the collected point cloud is 4X more sparse than KITTI's, of course, the amount of data is 4X decreased respectively. With the two datasets, we then fully evaluate the performance of the our proposed method for a total of 19 scenarios. Based on the KITTI testset, we choose four different sets of road driving test scenarios. And at the same time, in order to enrich the experimental content and verify our design effects, we conduct 15 experiments on Cooper using the T&J dataset.

Note that Cooper can also be applied to heterogeneous point clouds input. We elected not to conduct this test due to a lack of suitable LiDAR datasets.

We define single shot as point clouds collected by an individual vehicle, and cooperative sensing as merging all point clouds from nearby vehicles. We systematically analyze the test results of single shot and cooperative sensing to demonstrate the performance improvement on object detection. Qualitative results of cooperative perception under two experimental testsets are demonstrated in the following sections.

B. Evaluations on KITTI Dataset

In this section, we evaluate Cooper's performance using the KITTI dataset. As we know, KITTI provides raw consecutive 3D Velodyne point clouds in several scenarios. We choose one such segment sensing data in folder 2011/09/26/0009 as an example, shown in Fig. 2.

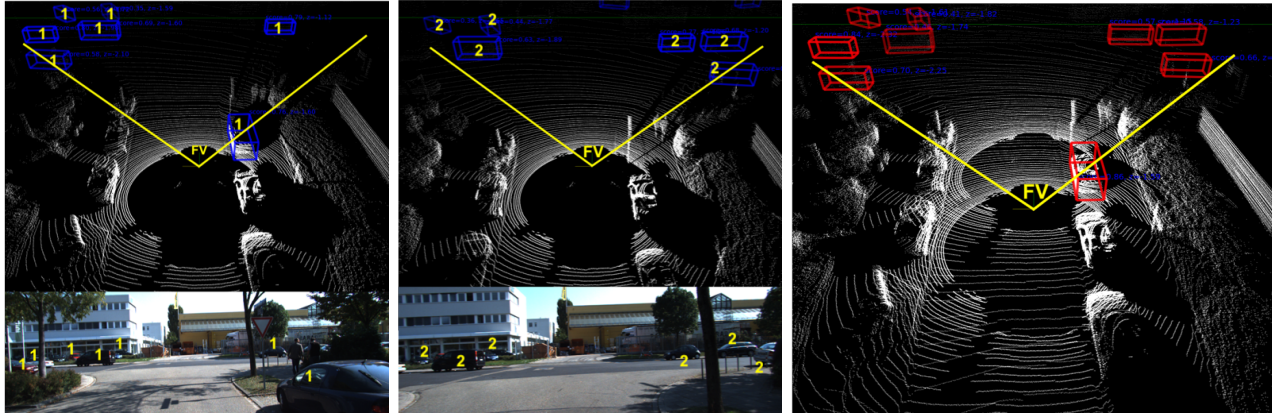


Fig. 2: Cooperative perception based on KITTI point clouds

To corresponding with 120° front view image, this LiDAR data of front-view area is evaluated. At beginning time t_1 , one single shot frame of 64-beam raw point cloud is collected in Fig. 2a. As the testing vehicle is moving forward after two seconds, another single shot frame of 64-beam raw point cloud is collected in time t_2 shown in Fig. 2b. By merging t_1 and t_2 's point clouds, we emulate the cooperative sensing process between two vehicles. We utilizes SPOD 3D detector to detect cars and draws results in red boxes to bound detected cars' location in Fig. 2c. Meanwhile, in order to compare the detection results on Cooper, we also adopt SPOD in Fig. 2a and Fig. 2b to detect cars, the results are drawn in blue boxes. We can find there are two improvements in this case study. First, we can see that in t_1 we observe 6 blue boxes, and in t_2 we observe 6 blue boxes yet again. However, when combined, we observe a total of 9 detected cars in merged data shown as red boxes, which includes all the results as blue boxes shows in t_1 and t_2 . This means the sensing area is extended by data sharing. The second is the detecting score/confidence value is increased on specific object. For example, on the right side, there is a nearby car is detected in t_1 , and in Fig. 2c, the same car is detected and the detecting score is increased by 10%. We also provide the corresponding images as ground truth in the bottom of Fig. 2a and Fig. 2b. The following is calculating the number of object detected by single shot and cooperative sensing, then, we compare against the ground truth in images for each case respectively. Cooperative perception process is evaluated on different distances in four scenarios: T-junction, stop sign, left turn and curve conditions in Fig. 3. Every three columns represents one collaboration process, which is similar to the example we demonstrated in Fig. 2. We draw the distribution of detection results using cells in each column. The number in each cell is the detecting score (0-1), the higher the score, the more positive the result. The symbol X represents a missing detection. The cell without score means

the object is out of detection area. Also there are different colors to indicate the distance. The darker the color, the farther the distance. According to the actual detection distance of LiDAR, we divide it into three scales of near ($<10\text{m}$), medium ($10\text{-}25\text{m}$) and far ($>25\text{m}$), which are represented by white, gray and black respectively in the illustration. We can find out that detected cars' quantity in cooperation is equal to or exceeds the number in both two merged single shots. Then,

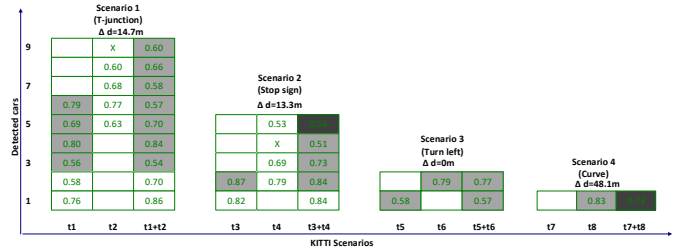


Fig. 3: Detection results distribution on KITTI

we use qualitative results to analyze the performance on the number and accuracy of detected vehicles shown in Fig. 4. The proposed cooperation method not only detects more cars, but also grants better detection accuracy because there is no missing detection in the cooperative point clouds.

C. T&J Dataset

Unfortunately, KITTI dataset does not provide enough experimental scenarios for proper Cooper testing because it is a vision benchmark collected by isolated instances of single vehicles. We are committed to multi-vehicle cooperation and sharing research, and thus, to improve the driving safety and experience of CAV, we focus on building a dataset that is suitable for collaboration, naming it the T&J dataset.

Our testing cars are equipped with high precision sensing systems, such as LiDAR system, radar system, vision system,

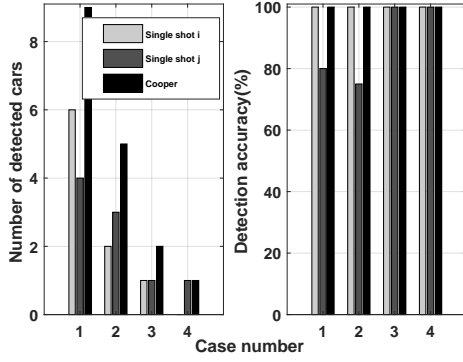


Fig. 4: KITTI detection results

and supplemental system such as GPS and IMU sensors. More specifically our sensor framework consists of the following sensors:-

- 2 X Front View Cameras
- 4 X Surround-view fish eye cameras
- 1 X inertial and GPS navigation system
- 1 X front-view 120° Rader
- 1 X Velodyne VLP-16 360° LiDaR
- 1 X Nvidia PX2

Velodyne VLP-16 360 LiDAR [3] is used for object detection and environmental mapping along with Radar, which utilizes radio waves to measure distance, and performs well in extreme weather conditions. However, LiDAR provides low resolution image information. Cameras, on the other hand, provides very high resolution image information, but, it fails to perform in extreme weather or environmental conditions. Four fish-eye lens cameras are used to perceive and navigate the surrounding environment. IMU sensors provides the system that monitors the dynamically changing movements of the vehicle. Also, GPS sensor data can be used to obtain a rough estimate of the location or the positioning of the car. Nvidia Drive PX2 [23] is a scalable AI supercomputer for our autonomous driving. This is a highly computational platform which combines deep learning, sensor fusion, and surround vision to change the driving experience. This can facilitate data fusion from multiple cameras, as well as LiDAR and radar sensors. Similarly, it uses Deep Neural Networks (DNN) [6] for the detection and classification of objects which dramatically increases the accuracy of the fused sensor data.

D. Evaluation on T&J Dataset

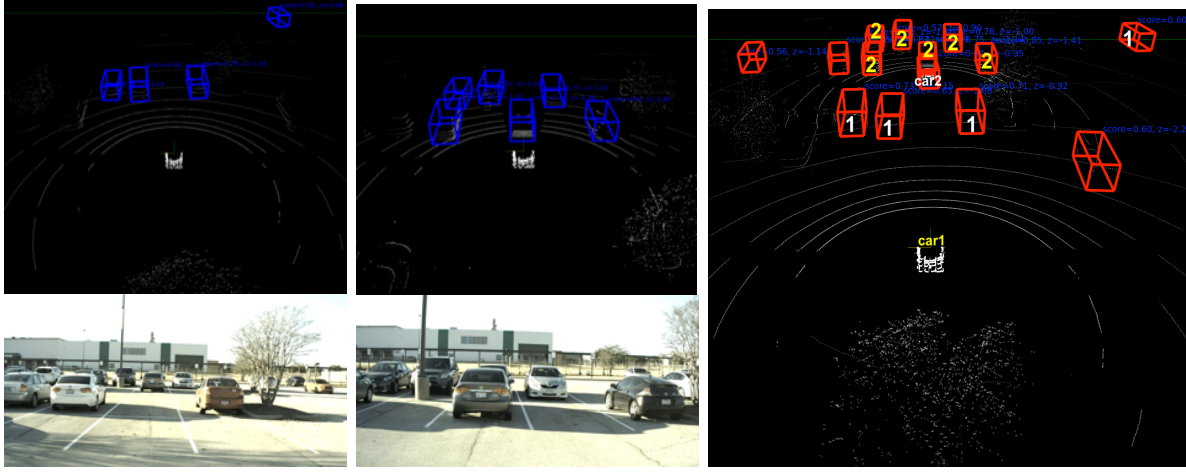
Next, we employ Cooper on our T&J dataset. We select a sequence of continuous frames of front-view perspective of point clouds as a example in Fig. 5. It can be clearly distinguished from the resolution of image that our point could is much more sparse when compared against the previous images. This is a car detection scenario in a parking lot. In one single shot, a frame of 16-beam raw point cloud is collected in Fig. 5a. In another single shot, another frame of raw point cloud is collected at close distance shown in

Fig. 5b. By merging these two frames of point clouds, we produce two vehicles' cooperative sensing. 3D detector detects cars and draws results in red boxes to bound detected cars' location in Fig. 5c. Similar to Fig. 2, SPOD detects cars in two single shots and draws in blue boxes respectively. SPOD also draws results in red boxes to bound detected cars' location in cooperative sensing. Meanwhile, ground truth images are shown in Fig. 5a and Fig. 5b. By studying this case, we can conclude that sensing area is expanded by data sharing because Fig. 5c detects all the objects exist in single shot. And most significant part here is that through Cooper, we see that the presence of new cars are discovered, cars that were not presence in the previous single shot. **This phenomenon is a direct proof to the shortcomings of fusion on a high level, object level fusion. Due to neither vehicles detecting the objects by themselves, there stands no possible way for the high level fusion to detect the objects that were missed. This, we avoid and overcome with low level fusion.**

We marked label 1 for the first shot's detected cars. Similarly, label 2 represents the detection result of the second shot. It is worth noting that there are three unmarked vehicles appearing in Fig. 5c. This is a significant discovery as this phenomenon indicates an increase in the detection capability of cooperative perception. We can extrapolate and assume that by receiving the perceptual information of nearby partners, CAV can greatly enhance its own range of perception, allowing for better detection of traffic information.

T&J dataset provided four sets of testing data, which were collected on the roads around our campus's parking lots. In these four scenarios, we conduct cooperative perception experiment. Different from KITTI test, in each experimental scenario, we sample the fusion data at different distances, so as to better display the disparity of information collected by vehicles in different regions. As Fig. 6 shows, in each scenario, we list detailed detection results of cooperative perception at different distances. Similar to Fig. 3, every three columns corresponding SPOD detection on two single shots and one cooperative sensing, represents a cooperative perception case. The test car can receive both nearby sensing data and relatively long-distance sensing data. For example, in Fig. 6a, from left to right, there are three cases in which a vehicle cooperates with other three located at three distance. It can be seen that in the cooperative perception of adjacent areas, such as the left cases in Scenarios 1 and 4, the individual detection results of two single shots are similar, but both output undetected targets, because these targets are blocked by unknown means in the single shots. Through cooperative perception, point clouds of blocked area are supplemented by each other, thereby these targets are detected. Moreover, the detected targets both shots, after cooperative perception, have a marked increased in detection scores. We evidence this phenomenon due to the redundancy of data and the presence of more features are gathered by harvesting detailed point clouds.

In all scenarios in Fig. 6, we carry out the cooperative perception of two cars, both are relatively far apart from each other. As a result, the detection area is expanded even larger.



(a) In one single shot, applying SPOD on 16-beam point clouds in front-view area to detect cars. (b) In second single shot, the detection results are drawn in blue boxes, bottom image is ground truth. (c) Cooperative sensing combine two single shots. The detected cars are drawn in red boxes using SPOD 3D detector.

Fig. 5: A cooperative perception example illustrated using T&J dataset

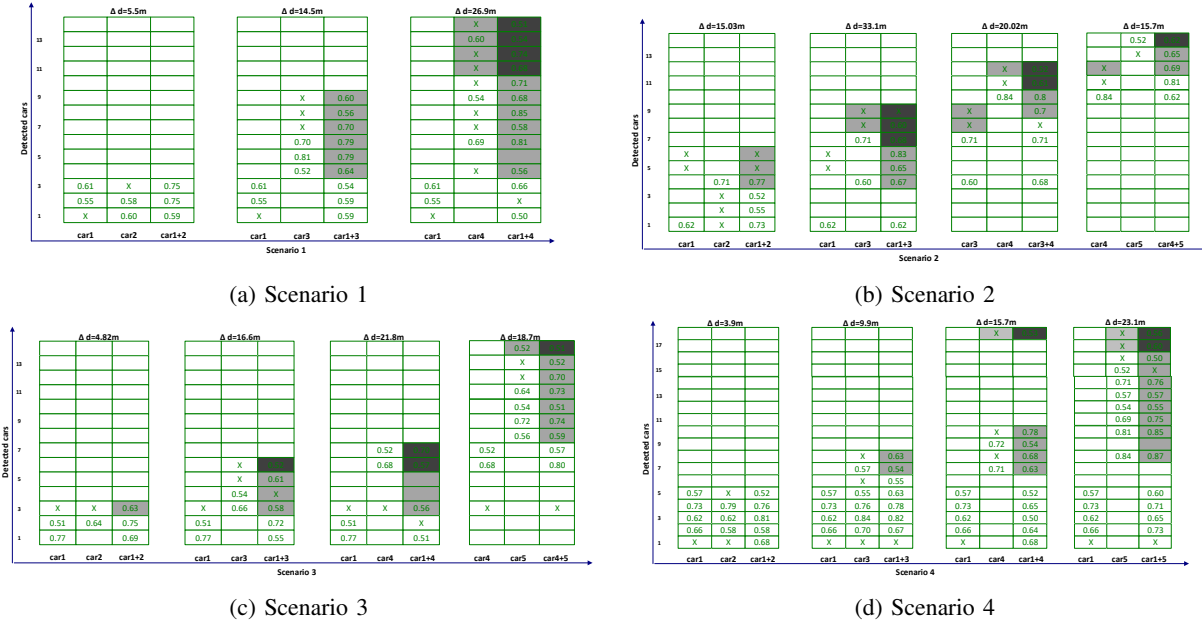


Fig. 6: Distribution of detection results on T&J testing scenarios

Every car can detect the target in front of itself. But for distant targets, they are powerless due to scarcity or blockage of point clouds. Cooperative perception enables global detection of objects located at far, medium, and near distance. Objects are appeared in cells of different colors. Similarly, some objects that are undetected by single shot are detected in cooperative sensing. This reinforces the fact that some objects that were not detected through traditional means can be discovered through data fusion. This shows that our design can complement some key features. This is a significant discovery on cooperative perception.

Then, we use qualitative results to analyze the performance

on the number and accuracy of detected vehicles shown in Fig. 7. From Scenario 1, we have the single shot analysis results for three different cases. It is clear that the number of cars detected based on the fused data is much higher than either of the cars alone. Despite the high detection rate however, we do see that even while fused there are still some cars not being detected as shown by the corresponding chart.

In Scenario 2, we find that there is a high amount of cars that is hard to detect from either car alone, but shows up when fused. This change of environment hold high relevance to common place areas such as a full parking lot or congested junctions where each car is limited by the cars around it.

Should there be a speeding car that is ignoring stop signs or running the red light, the fusion will mitigate the likelihood of a missed detection for all cars involved in the immediate vicinity.

In both Scenarios 3 and 4, we find that, similar to the trend shown in Scenarios 1 and 2, we have a closely related relationship between fusion and increase in object detection. As each scenario takes place in different environments, time of the day, different levels of congestion, the fusion method is proven robust and is able to adapt to different environments while retaining its capabilities to augment the status quo.

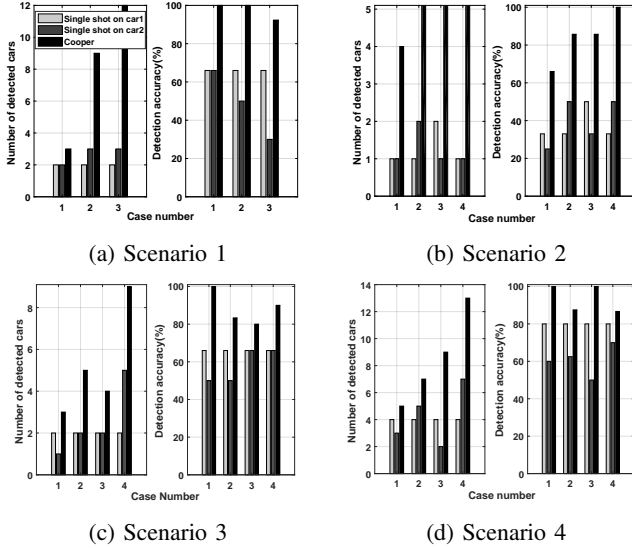


Fig. 7: Analysis on detection results in T&J testing scenarios

E. Statistical Analysis

Our statistical test results show that in the experimental scenarios of KITTI and T&J datasets, some of the targets in cooperative perception are detected by both, some by only one, and some are detected by neither. Detection difficulty is classified as easy, moderate and hard respectively. Specifically, easy refers to when one or more vehicles are able to detect the same object. Moderate refers to when only one vehicle is able to clearly detect this object. Finally, hard is given when no cars are able to detect this object.

In Fig. 8, we calculate the improvement of detection performance on these three types of objects. For example, from the line marked easy, we see that we have an improvement of 10% in detection score 80% of the time. Taking the direct implication of our testing, we see that the detection scores for easy and moderate achieve a marginal yet consistent increase in detection rate; mainly distributed within 10% in detection score improvement. This is because both easy and moderate objects contain detailed and saturated sized point clouds captured from a single scene, resulting in the fusion providing only marginal improvements to the detection results.

However, note that when we test the third type of object, the hard object detected by neither, we find that we are consistent

with our findings that we have above, **our detection score improvement is a flat increase of 0.5 in raw detection score at worst and just this alone is enough for autonomous vehicles to note the object for avoidance prevention, because they only need to know that there is an object there where previously one was not discovered.**

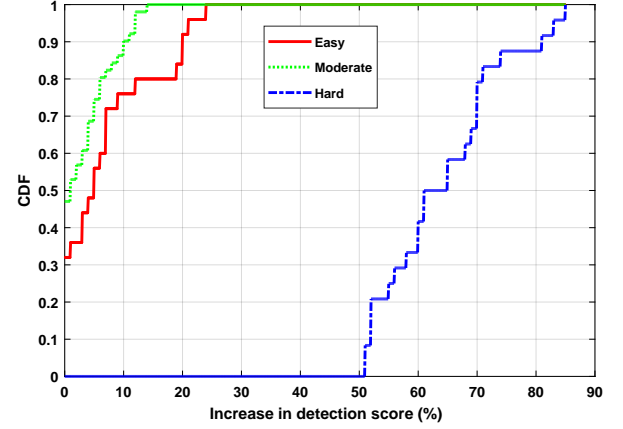


Fig. 8: Detection performance improvement by cooperative perception

We record time cost of detection based on single shot and cooperative data, shown in Fig. 9.

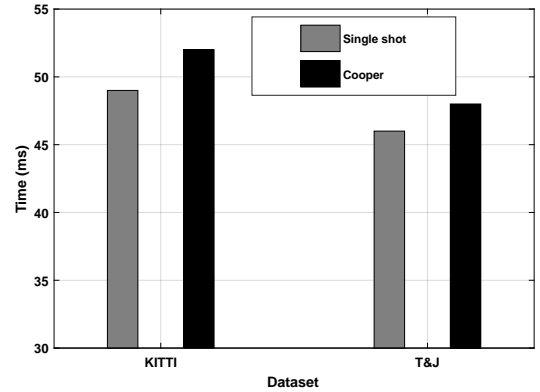


Fig. 9: Detection time cost on single shot and cooperative sensing in our experiments

As latency impacts the performance and reliability of all autonomous vehicles heavily, we tested our fusion method against both the KITTI data and our own. The SPOD model for 3D car detection is executed on the GeForce GTX 1080 Ti GPU [1]. In both experiments, we compared the time cost of object detection for both single shots against the fused data. **In both instances, fusing the data used 5 ms over the baseline data, a very minimal increase in detection time for a significant increase in the number of objects detected.**

F. Fusion Robustness

From a realistic standpoint, we will inevitably have to deal with sensor drift, so to deal with this phenomenon, we must test our fusion method of robustness against sensor drift. When

integrating GPS and IMU, we observe yields of less than 10 cm in positional errors [5]. To test the robustness of our fusion method, we conducted procedural artificial skewing of our GPS readings. We skew the GPS data as follows:

- Skewing both x and y coordinates to the maximum bounds of known GPS drifting.
- Skewing just one axis to the limit of GPS drifting.
- And pushing past that boundary by doubling the maximum GPS drifting to simulate abnormal instances.

With the GPS readings skewed, we then tested the detection score for each of the different type of drifting scenarios against the baseline GPS reading. As evinced from Fig.10, we see that with the exception of already known undetected vehicles, we have a similar clustering of the skewed detection scores versus the baseline score, with the overwhelming majority achieving successful detection.

It should be noted, however, that skewing the readings surprisingly improved the detection score in several instances, possibly masking the inherent drift from the baseline GPS reading. And just as some of the skewing helped the result, it also caused the detection to fail for two instances.

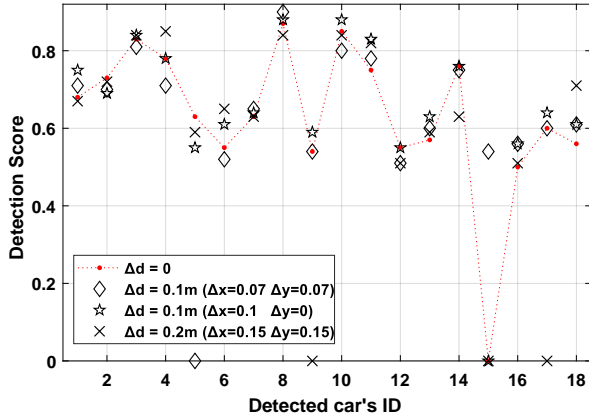


Fig. 10: Cooperative perception results on GPS reading drifting

G. Networking Requirement

Even though point clouds can be simplified to coordinate values, we still need to consider the gap between data generated by autonomous vehicles and the limited wireless networking throughput, such as the limited bandwidth provided by DSRC [11]. We adopt a strategy to extract data from region of interest (ROI), such as traffic lights, blocked areas, nearby vehicles and free-space in driving path, to further reduce data size to hundreds *KB* per frame. Background data like buildings, trees are subtract. Because these information can be constructed by each vehicle after several times mapping measurement. This allows for retention of valuable information of immobile objects while keeping the size of the ROI data small. For object detection purpose, ROI data will be extracted whenever failure detection happened on this area.

However, just knowing the relative ROI is not optimized enough. The ideal case is to have a multitude of real world ROI categories that provide a guideline for the bases of how much data is needed for an optimal balance of data size versus detection accuracy. To illustrate the importance of this tradeoff, we present three different types of ROI categories and their respective data consumption via Fig.11 and Fig.12 respectively where the sample rate in the latter is 1Hz, or 1 frame per second; we simulated and gathered the total data consumption between two cars, both utilizing a 16-beam LiDAR, every second over an eight second time frame. Note, we observe that message exchange rate for cooperative perception does not require as high a sensing rate as the standard rate for individual vehicles. Because for easy or moderate objects, detailed sized point clouds are already captured. While due to blocking or distance, we may experience an insufficiency of point clouds, making objects hard to detect. In most cases, the native data on a recipient vehicle only needs to be supplemented by a single data frame from different view perspective. Excessive exchanging of frequencies only leads to unnecessary data, hence needlessly congesting the communication channels. With efficiency and lightweight traffic as a constraint, we decided that a sample rate of 1 frame per second is enough to satisfy the needs of Cooper whilst remaining within our set of constraints.

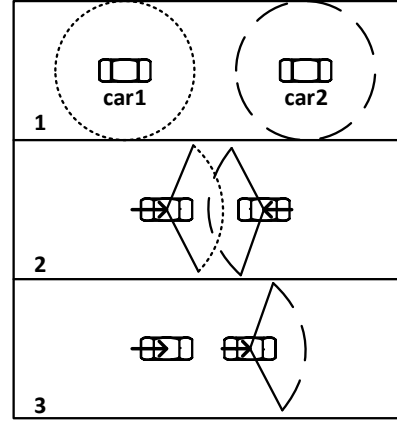


Fig. 11: Region of interest

As seen in Fig.11, we have three different scenarios, each representing a general phenomenon. For the first one, we see that two cars are fairly apart from each other, laterally but fairly close horizontally. We would typically see this situation in two lane drive with opposite directions separated by a single lane divider. In this scenario, we would ideally want as much information as possible as we lack the safety of a physical buffer between the vehicles. In situations like this, we transfer the entirety of the frame of LiDAR data and this is the most costly of all scenarios as evinced by Fig.12. From the same scenario, we can calculate that for the most expensive data transaction, the total data size can be compress into around

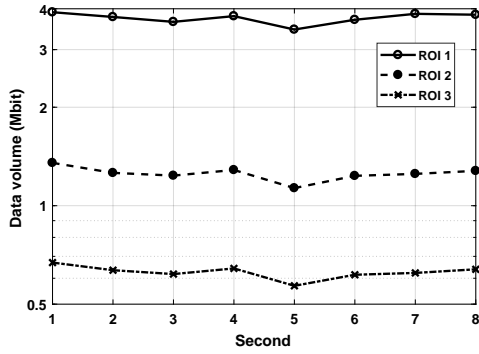


Fig. 12: Data volume exchanged between two cars

1.8 *Bbit* per frame for each car.

Next, we have the case of cars in closer lateral proximity to each other, representing typical junctions where cars from all directions are able to see the opposing car. In situations such as this, the ROI is typically the field of view from the driver's perspective, making only a 120 degree field of view our minimal requirement. As both vehicles need to exchange this information, the transaction cost is additive for each of the participating vehicles.

Lastly, we have the most common situation of one car needing the field of view of a leading car. The trailing car is the one needing the information and thus the transaction is one way, consuming the least amount of bandwidth out of all three scenarios.

Thus, deriving from the simulation of the three different cases, the three presented are within the capacity of DSRC bandwidth.

In summary, we prove that Copper method outperforms individual perception on extending sensing area, improving detection accuracy and complementation of object detection. We find that collaboration offers more information, even some are not perceived by individuals. The most important, we demonstrate that the bandwidth of DSRC can satisfy point clouds transmission for cooperative perception. We would like to mention that our design succeeds in privacy preservation because only LiDAR data is involved for sharing.

V. RELATED WORK

Rapid development of autonomous vehicles has motivated research institutions to develop representations to perceive local environment, such as lane detection, traffic sign detection and detect objects like cars, cyclists and pedestrians [18], [19], [29], [17] based on the open datasets [8], [7], [9]. As we know, these datasets are collected by multiple sensing devices from individual vehicles. To achieving self-driving, we put heavy emphasis on accuracy cognition of the surrounding local environment. However, the detection results still have vast room for improvement even when utilizing state-of-the-art Convolutional Neural networks (CNNs) [12].

Current works make use of low level fusion of sensors to extract the features or objects for purpose of tracking [13]. However, this does not incorporate the use of raw data as is for the purpose of fusion and object detection. Papers such as [10] and [4] discuss methods of fusion that constructs theoretical architecture for low level fusion and detection.

To the best of our knowledge, there are no prior work done to implement the concept of multi vehicular raw sensor data for the purpose of object detection.

This room for improvement is also the cause of severe consequences because self-driving cars may make wrong decisions due to failure of detection of objects. A Cooper framework for connected autonomous vehicles can solve the aforementioned issues through cooperative sensing. However, none of the public datasets and related detection methods explicitly consider low level fusion approach as a solution.

VI. CONCLUSION

We propose Cooper for connected autonomous vehicles as an entry to a broader platform for CAV. This method facilitates a CAV capable vehicle to combine its sensing data with that of its cooperators to enhance perceptive ability, thereby improving detection accuracy and driving safety. In order to reconstruct local environment, we map point clouds into their corresponding object positions. This will merge and align the shared data that is collected from nearby vehicles, which may provide data scopes coming from different positions and angles. We incorporate deep learning based SPOD with Cooper to detect 3D objects from aligned LiDAR data, marking and discovering previously undetected objects. Finally, we evaluated Cooper on KITTI and our collected dataset, showing that the Cooper is capable of enhancing detection performance by expanding the effective sensing area, capturing critical information in multiple scenarios and improving detection accuracy.

REFERENCES

- [1] Geforce – nvidia. <https://www.nvidia.com/>.
- [2] Intel – the coming flood of data in autonomous vehicles. <https://www.intel.com/content/www/us/en/automotive/autonomous-vehicles.html>.
- [3] Velodyne – lidar. <https://velodynelidar.com/>.
- [4] M. Aeberhard and N. Kaempchen. High-level sensor data fusion architecture for vehicle surround environment perception. In *Proc. 8th Int. Workshop Intell. Transp.*, 2011.
- [5] K.-W. Chiang, T. T. Duong, and J.-K. Liao. The performance analysis of a real-time integrated ins/gps vehicle navigation system with abnormal gps measurement elimination. *Sensors*, 13(8):10599–10622, 2013.
- [6] D. Cireřan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. *arXiv preprint arXiv:1202.2745*, 2012.
- [7] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.
- [8] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [9] X. Huang, X. Cheng, Q. Geng, B. Cao, D. Zhou, P. Wang, Y. Lin, and R. Yang. The apollo-scapes dataset for autonomous driving. *arXiv preprint arXiv:1803.06184*, 2018.

- [10] N. Kaempchen, M. Buehler, and K. Dietmayer. Feature-level fusion for free-form object tracking using laserscanner and video. In *Intelligent vehicles symposium, 2005. Proceedings. IEEE*, pages 453–458. IEEE, 2005.
- [11] J. B. Kenney. Dedicated short-range communications (dsrc) standards in the united states. *Proceedings of the IEEE*, 99(7):1162–1182, 2011.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [13] R. Labayrade, C. Royere, D. Gruyer, and D. Aubert. Cooperative fusion for multi-obstacles detection with use of stereovision and laser scanner. *Autonomous Robots*, 19(2):117–140, 2005.
- [14] B. Liu, M. Wang, H. Foroosh, M. Tappen, and M. Pensky. Sparse convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 806–814, 2015.
- [15] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [16] National Transportation Safety Board. Preliminary report. <https://www.nts.gov/investigations/AccidentReports/Reports/HWY18M-H010-prelim.pdf>, 2018.
- [17] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas. Frustum PointNets for 3D Object Detection from RGB-D Data. *ArXiv e-prints*, Nov. 2017.
- [18] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [19] J. Ren, X. Chen, J. Liu, W. Sun, J. Pang, Q. Yan, Y.-W. Tai, and L. Xu. Accurate Single Stage Detector Using Recurrent Rolling Convolution. *ArXiv e-prints*, Apr. 2017.
- [20] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [21] B. Schwarz. Lidar: Mapping the world in 3d. *Nature Photonics*, 4(7):429, 2010.
- [22] J. Shi, W. Wang, X. Wang, H. Sun, X. Lan, J. Xin, and N. Zheng. Leveraging spatio-temporal evidence and independent vision channel to improve multi-sensor fusion for vehicle environmental perception. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 591–596. IEEE, 2018.
- [23] Wikipedia. Nvidia – wikipedia, the free encyclopedia. <https://en.wikipedia.org/wiki/DrivePX-series/>, 2016.
- [24] Wikipedia contributors. Rotation matrix — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=Rotation_matrix&oldid=875545324, 2018.
- [25] Wired. Why tesla’s autopilot can’t see a stopped firetruck. <https://www.wired.com/story/tesla-autopilot-why-crash-radar/>, 2018.
- [26] B. Wu, A. Wan, X. Yue, and K. Keutzer. Squeezeseg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from 3d lidar point cloud. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1887–1893. IEEE, 2018.
- [27] Y. Yan, Y. Mao, and B. Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018.
- [28] Q. Zhang, Y. Wang, X. Zhang, L. Liu, X. Wu, W. Shi, and H. Zhong. Openvdap: An open vehicular data analytics platform for cavs. In *Distributed Computing Systems (ICDCS), 2017 IEEE 38th International Conference on. IEEE*, 2018.
- [29] Y. Zhou and O. Tuzel. VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection. *ArXiv e-prints*, Nov. 2017.