

Is the Synthesized Scene in Autonomous Driving Realistic?

Elaine Yao

University of British Columbia

Abstract

For Autonomous Vehicles(AVs), sensor perception is safety-critical. Failures in object detection can cause disasters. Despite various prior works on adversarial 3D physical attacks in AVs, all of them are simulating placing obstacles on the road with a synthesized scene. However, the rendering functions to generate the synthesized scenes may not be able to ensure the physical consistency between the obstacle and the background. The real-world sensor perception is much more complicated. In this project, we present the study of authenticity issues of rendering-based synthesized scenes in AV systems. We evaluate this by generating different synthesized scenes and testing the performance of the neural network on them. This allows us to quantify how realistic these scenes are and understand the impact of physical consistency in the scene on the effectiveness of generated adversarial objects.

We design a comprehensive test suite aiming at evaluating whether the method to integrate a 3D object in the road background is realistic or not. We adopt empirical approaches that address four main design challenges: various impact factors, physical environment consistency, domain-specific metrics, and automated pipeline. We evaluate the synthesized scene with our test suite in representative open-source industry-grade AD system object detection models with real-world driving scenarios. We also choose a state-of-art adversarial 3D physical attack for evaluation in malicious cases. Our results show that most synthesized scenes are not realistic enough so the object detection fails to detect the obstacles in it. Such a phenomenon can reduce the effectiveness guarantee of generated 3D adversarial attacks in the physical world.

1 Introduction

Autonomous Vehicles(AVs) can sense the surrounding environment and move safely with little human input. They are

playing an important role in future transportation. Large companies such as Google, Uber [6] are racing to develop AVs and some high-level, such as level 4 self-driving cars have already been deployed on the road. Level 4 is considered to be fully autonomous driving. It can handle complex urban driving situations without driver intervention. A fundamental part of the autonomous driving system is perception. It uses sensors [9] such as cameras, LiDARs, Radars, IMU(Inertial Measurement Unit) and GPS to know the physical environment and react accordingly. Among them, perception sensors including cameras and LiDARs provide the obstacle and traffic sign information to AVs to avoid wrong decisions like collisions and violating traffic rules, etc. Failures in perception can pose a threat to the safety of self-driving. In 2020, a tesla car in autopilot mode collided with an overturned truck as it failed to detect it. Therefore, multiple prior works have been studying the security of these perception sensors.

Prior work has shown that AVs are vulnerable to attacks on camera [13, 17, 33] or LiDAR sensors [10, 12, 28, 35]. Adversaries can change the texture of a 2D image [33](e.g., stop sign) or add well-designed adversarial patches [17] to mislead the cameras. They can also inject laser [12] to spoof the LiDAR sensors.

All of these studies, however, are limited to attacks with synthesized background, i.e., integrating the adversarial object with the road background through rendering functions instead of realistic simulation [11, 26, 28, 33, 35]. To simulate the scenario where a 3D vehicle is put on the road, This work will synthesize the attacked-influenced sensor perception, for example, the point clouds by LiDAR and images by a camera with respective 3D rendering functions. These rendering techniques provided by computer graphics can simulate the real sensor functions but still lack comprehensive consideration in sensing the environment due to the simplicity of its model. By contrast, sensor perception in the physical world can integrate more information such as light conditions, the realistic texture of adversarial objects, and reasonable positions. Although most recent works in adversarial attacks will evaluate their methods in the physical world to show the effectiveness

and feasibility, they use the synthesized method to generate the malicious objects at a faster speed. The assumption that the synthesized attack-influenced background is physically consistent should be believed to hold in general [11] and thus examining the effectiveness of this method is an urgent call.

This project presents a study on the evaluation of the reality of synthesized backgrounds in AD perception systems today. We test the above rendering-based simulation assumption by evaluating the neural network performance on these integrated sensor perception outputs used in the state-of-art adversarial attack work. This allows us to gain a solid understanding of how much authenticity guarantees the use of synthesized background can provide a realistic simulation way to generate effective adversarial objects. Specifically, we consider physical 3D objects as the attack vectors for real-life feasibility and examine the performance of object detection neural network models deployed in real AV systems on the synthesized scenes.

Even though previous works have designed perception rendering functions for cameras and LiDAR, we find that simply feeding them with different objects and backgrounds won't meet our requirements. First, we need to identify the factors in the synthesized scene that might influence the detection accuracy of the neural network. For example, the object detection model may find it difficult to detect an object which is far from away the sensor. Also, the color of the obstacle is similar to the surrounding environment so that it's hidden from the neural network. Second, physical consistency between the obstacle and the driving background should be maximum guaranteed. No matter where we put the obstacle, it should stand on the road instead of floating in the air or hitting the ground. It should also follow the shadow caused by sunlight. Third, to quantify the authenticity of the synthesized scene, we need to come up with domain-specific metrics. The previous works use different metrics to measure whether their adversarial obstacle achieves the goal and lacks a unified standard, which makes it difficult to provide fair and reasonable metrics. Fourth, we need to develop an automated pipeline for generating different synthesized scenes, evaluating the neural network performance in the scene without attack as well as with attack. Manually adjusting the parameters can take a long time and it's hard to do large-scale analysis.

Towards this end, we design an automatic and comprehensive synthesized scene test suite, which addresses the challenges above and thus provides evaluation for the authenticity of these rendering methods. Through preliminary experiments, we choose different driving backgrounds, 3D obstacle properties(including the color, shape, and texture), and the interaction between the background and the obstacle, e.g., the relative position, as the impact factors and serve them as the parameters to adjust. The attackers assumed in the previous work can just place an object on the road as simulated in the synthesized scene. To systematically generate a realistic scene, we adopt camera imaging theory to adjust the height

of the object so that it's standing on the road. Light condition is considered to comply with the driving background. Also, we start with a normal obstacle that can be obtained from life easily, e.g., a common chair. Under these test settings, we address design challenge 3 by considering the correctness of the bounding box of a detected object, object class, and its corresponding confidence score. We extract these by parsing the output of the object detection neural network. Also, we use these as building blocks to compute the overall scores for the authenticity of the scene. In the end, we developed automated pipelines for selecting different factors and evaluating the detection performance under benign and malicious cases.

We evaluate the scene synthesizing method in MSF-ADV [11] and choose the image object detection, neural network model, in Autoware.AI [3], which is representative of current AD systems. We also choose the attack in MSF-ADV [11] to generate adversarial 3D objects which can both fool the camera and LiDAR object, detection models. We select 3 shapes of chairs from McGill 3D Shape Benchmark [7] and evaluate each on 5 real-world driving scenarios from the KITTI dataset [21]. 60 different scenes are synthesized and evaluated. Our results show that the benign obstacle in the synthesized scene fails to be detected in all the test settings. We also find that for the attack strategy generating the adversarial object, if the benign object fails to be detected in the first place, it's also hard to generate effective adversarial obstacles. What's more, in this situation, it's hard to provide a guarantee that the generated adversarial object is effective in the physical world.

In summary, this work makes the following contributions:

1. We study the authenticity of synthesized scenes in AD perception systems. We successfully design a comprehensive test suite aiming at evaluating whether the method of integrating a 3D object in the road background is realistic or not.
2. We adopt empirical approaches that address four main design challenges: various impact factors, physical environment consistency, domain-specific metrics, and automated pipeline.
3. We evaluate the synthesized scene with our test suite in representative open-source industry-grade AD system object detection models with real-world driving scenarios. We also choose a state-of-art adversarial 3D physical attack for evaluation in malicious cases. Our results show that most synthesized scenes are not realistic enough so the object detection fails to detect the obstacles in it. Such a phenomenon can reduce the effectiveness guarantee of generated 3D adversarial attacks in the physical world.

While rendering the obstacle into the road background is a general way of generating adversarial 3D obstacles, prior works lack the realistic validation of the synthesized scene.

In this project, we try to evaluate it by measuring the performance of the neural network under different settings. We hope that our findings can inspire more future related research to validate their rendering process in AD perception when designing the 3D adversarial obstacles. The GitHub repository is publicly available at <https://github.com/ElaineYao/571p>

2 Background

2.1 AV Perception System

In the state-of-art AV systems, perception plays an important role in ensuring safety. The sensors are used to detect the obstacles and measure the velocity or distance in real-time. Typical systems in high-level, such as Level 4 [1] AVs adopt both LiDAR and camera for visual perception. LiDAR [15] can detect the ranges by shooting an object with a laser and measure the distance by getting the time for the light to be reflected by the receiver. Compared with RADAR, LiDAR is much more accurate in resolution. Thus it's used to reconstruct exact 3D models of objects in autonomous systems. However, it's difficult to get texture-related information such as the color [19]. On the other hand, camera images are good at providing shape and texture information but lack depth and distance information due to their 2D imaging. To compensate for the weaknesses and utilize the strength of each sensor, most AV systems will adopt Multi-Sensor Fusion(MSF) design, in which it will fuse the sensor reading from both LiDAR and camera.

Figure 1 shows an overview of the perception module in a common AV system [4]. 3D objects are first perceived by LiDAR and camera to generate point clouds(LiDAR) and frames of images(camera). These raw sensor outputs will then go through a pre-processing unit for the aggregated feature and ROI(Region of Interest) extraction. Pre-processed features will be fed into the LiDAR perception network and camera perception network respectively to get the detection results. The MSF algorithm will fuse the outputs of two perception networks and give the final detection output.

In this project, we focus on the camera perception part. 3D objects are sensed by the camera in the form of 2D images. When a random obstacle is put in the middle of the road, a synthesized image with the obstacle and road background is generated. This image will then be fed into the object detection neural network for results.

2.2 Synthesized Scene

To synthesize the obstacle with the road background, many prior works [11] use Neural 3D Mesh Renderer(NMR) for camera rendering [22]. NMR provides a way to generate a 2D image from the 3D world. It can transfer rendering gradient with consideration of texture, lighting, camera, and object

shapes. Given the camera pose, light condition, relative position between the camera and the overall background, NMR will provide the image output of this background with a certain camera setting.

Fig. 2 overviews the image synthesizing process in MSF-ADV [11]. It first chooses the background from the target road and the obstacle that it intends to put on the road, for example, a brown chair. The 3D chair is presented in the form of the point cloud, which is a set of data points in space. Camera parameters and light conditions will be set to render this chair into a 2D image. Then the relative location of the chair in the background is set. Original pixels in the background will be masked by the pixels in the 2D chair image to simulate how a chair is put in the middle of the road. Before feeding the synthesized scene into detection neural networks, some pre-processing steps such as data transformation, utilizing the Region of Interest(RoI) filter to clear unrelated input parts and collecting aggregated features are performed. These pre-processing processes can reduce the input size fed into the neural network and greatly improve the inference speed.

In this project, we will use this as the target pipeline for generating synthesized scenes and will evaluate the authenticity of the synthesized output.

2.3 Adversarial 3D Object Attacks

Prior works find that it's easy to fool models with deceptive data, which causes the malfunction in the neural network model. This kind of adversarial attack is also studied in the context of the physical world [16, 18, 32]. In the AV systems, some prior works proposed physical adversarial attacks by placing obstacles in the air or on the road [16, 18, 32]. Some are targeting fooling the camera-based perception neural network [16, 18, 32] while others are focusing on LiDAR-based perception neural network [12, 27]. Recently, considering that MSF design is widely used in AD systems, MSF-ADV [11] is proposed to attack both LiDAR and camera sensors with a 3D adversarial object.

MSF-ADV [11] treats the process of generating the adversarial attacks as an optimization problem and Fig. 3 provides an overview of the process to generate adversarial obstacles. It first picks a normal 3D object and applies 3D transformations including rotation, and position shifting to get different angles of the object. This is to improve the robustness of the obstacle in various kinds of environments. Then it will generate a point cloud and image of this 3D object through ray-tracing [8] and NMR [22] to simulate the perception output of LiDAR and camera sensors. These two sensor outputs will be integrated with the road background and pre-processed before being sent to perspective perception neural networks and the MSF algorithm. The attacker is assumed to be able to perturb the shape and position of the 3D object. Also, the adversarial loss function is designed to cause the malicious object not to be detected by the MSF algorithm as well as keep the obstacle

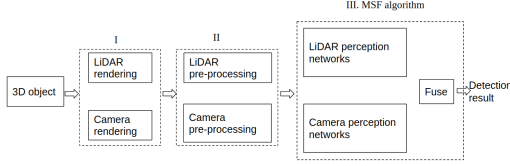


Figure 1: Perception module in AV systems

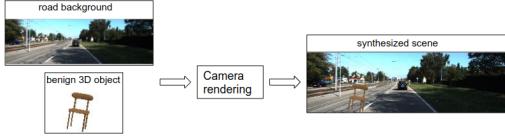


Figure 2: Synthesized scene through camera rendering

stealthy. After obtaining the generated malicious object, the attacker can just 3D print it and place it on the road according to the parameters in the optimization process.

In this project, we will use MSF-ADV [11] as the target adversarial attack in the evaluation part. MSF-ADV [11] uses the same method in the previous section to generate a synthesized scene. Therefore, we want to see whether the authenticity of the synthesized scene will influence the effectiveness of the generated malicious object.

3 Attack Model and Design Challenges

3.1 Attack Goal and Threat Model

Attack goal: Fail both LiDAR and camera perception neural networks in an AV system. In this project, we target at evaluating the effectiveness of the MSF-ADV [11] attack on the synthesized scene. The attack goal is straight-forward for affecting the safety of autonomous driving: fail the AV perception with LiDAR and camera sensors in the target AV, so that it can't detect the obstacle in front of it and collide with it. This attack also assumes the AV system has a fail-safe mechanism for an emergency brake. Even though the system is equipped with an Automatic Emergency Brake(AEB), prior works still show that the victim's vehicle can be hit by the cars behind if they fail to brake in time. Therefore, MSF-ADV

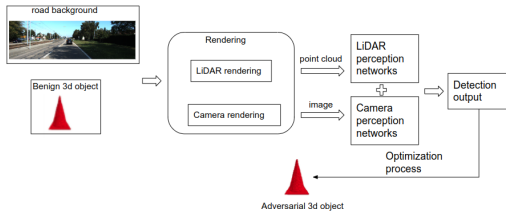


Figure 3: Pipeline for generating adversarial 3D object attacks for both Lidar and Camera

is designed for real-world physical attacks with the current industry-level AD systems.

AV systems that use more than 2 perception sensors will usually be equipped with an MSF algorithm, which is to fuse the information from all of the sensors and then make decisions. MSF still has the chance to correct the wrong sensor perception as long as there is at least one clean sensor source, which is unattacked. Since this attack aims at defeating the MSF-based AD perception, it must fail all the sensor perceptions for high attack effectiveness. Thus MSF needs to attack all visual sensors (i.e., camera and LiDAR) at the same time.

Threat Model. MSF-ADV [11] is designed in a white-box setting. It assumes that the attacker knows the MSF algorithm and the corresponding neural network model for each sensor perception in the target AV system. Many prior works that study physical adversarial attacks towards sensor perception in AV systems also have similar assumptions [18, 32]. This assumption is valid as the attacker can reverse engineer the purchased or rent AV perception module to get the algorithm details. What's more, lots of open-source industry-level AV algorithms can also be the victim, e.g., Baidu Apollo [4], Autoware.AI [3].

Since this attack is closely related to the road condition and background, the attacker is also assumed to take camera images and LiDAR data about the target road background to generate the adversarial objects using MSF-ADV. After getting the adversarial object in the target background, the attacker can 3D print the object and place it at the calculated place. To make the attack more powerful, the attacker can even fill it with high-density material like granite to make it heavier. As a result, when the car fails to detect this object and collides with it, the damage might be more severe. Fig. ?? shows the benign object, i.e., a chair, and the malicious chair generated by MSF-ADV [11]. We can notice that the surface of the malicious chair is more glitchy and it's easy for the vehicle ignores the overall shape and crash into it.

3.2 Design Challenges

In the state-of-art work in generating adversarial physical 3D objects, they all use rendering functions to simulate the sensor outputs and integrate the rendering output with the road background to get the attack-influenced scenarios. Although many works provide rendering pipelines, we find that simply feeding them with different objects and backgrounds won't meet our requirements due to 4 unique challenges:

C1. Need to identify factors in the synthesized scene which might influence the detection accuracy of the neural network. To evaluate the authenticity of the synthesized scene, we need to find the factors that might make this scene more realistic or the opposite. However, these factors are not fully explored in prior work. When an adversarial object is generated in a previous study, they only consider certain conditions, e.g., specific background, and fixed relative position.

For example, for 3d physical attacks, previous work predominantly considers the relative position between the obstacle and the target vehicle [35], the condition of target road background [11, 27]. Although there are many other factors such as the aspect ratio of the background and light conditions, they aren't taken into consideration in previous works when synthesizing the scene. One possible solution is to pick different obstacles and road backgrounds as many as possible. However, this not only adds up the testing overhead and thus decreases the testing efficiency, but also results in redundancy in testing cases, which might influence the trust in the final results. Thus, it is highly desired to identify some factors that can represent the difference between various synthesized scenes.

C2. Physical consistency between the obstacle and the driving background should be maximum guaranteed. To ensure the synthesized scene is generated realistically in the first place, prior works are trying to select a physically realizable object such as the drones holding a cardboard with reflective surface [35] and printable traffic cone [11]. Since these adversarial objects usually need a lot of optimization iterations to be generated, it's impractical to drive the vehicle and put the obstacle in real-life to get real-time adversarial outputs. Therefore, most prior work will synthesize the impacts of adversarial obstacles from the physical sensor to both image and point cloud outputs, as well as integrate them into the road background. Despite the consideration of the physical realizability of the object, we also need to model the physical consistency when integrating the object and the background. For example, no matter where we put the obstacle, it should stand on the road instead of floating in the air or even hitting the ground. This means we can't put the obstacle in the background with a random position, instead, the physical model between the obstacle size and relative position should be established. It should also follow the shadow caused by sunlight. For example, the road background is taken when there is light from the east, thus the original object in the background has the shadow facing west. When the obstacle is put on the road, it's also supposed to have a shadow in the same direction. Or the neural network may fail to recognize it because of the physical inconsistency.

C3. Domain-specific metrics for quantifying the authenticity of the synthesized scene. In previous works designing adversarial obstacles, they adopt the optimization-based method, an optimization loss function is designed to measure the effectiveness of the adversarial object, as well as the stealthiness of the malicious obstacle. These loss functions are designed based on the confidence value of the object detection neural network, which reflects the probability that the region contains an object [11]. Also, they use metrics to measure the smoothness of the surface of the object. However, these metrics are designed in a way to make the optimization process more convenient, and faster. Our preliminary experiments show that even after the optimization, it's still possible


$$\text{Intersection over Union (IoU)} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


Figure 4: Intersection over Union(IoU)

for the neural network to detect the obstacle. Therefore, we have to come up with new metrics to measure the performance of the object detection neural network and the adversarial attacks. Meanwhile, previous works use different metrics to measure whether their adversarial obstacle achieves the goal and lack a unified standard, which makes it difficult to provide fair and reasonable metrics.

C4. Develop an automated pipeline for generating different synthesized scenes. After obtaining the factors that influence the authenticity of the synthesized scenes, corresponding value ranges should be determined. Usually, the prior works only consider limited factors with certain value ranges, the rendering pipeline they're using may not be suitable for comprehensive testing. It's also unrealistic to manually choose the combinations among factors as this takes a lot of effort. Automated pipelines generating different testing scenarios should be developed. Meanwhile, since different synthesized scenes are generated to measure its authenticity through evaluating the neural network's performance in benign and malicious cases, this pipeline should also be equipped with neural network model inference and automatically compute the evaluation metrics along the way. Having this automated pipeline can help us to do a large-scale analysis on how realistic these synthesized scenes are.

4 Approach

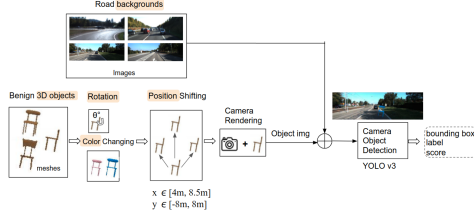
In this project, we address the 4 aforementioned challenges by designing an automatic and comprehensive test suite for synthesized scenes, which provides ways for evaluating the authenticity of the rendering methods used in MSF-ADV [11].

4.1 Overview

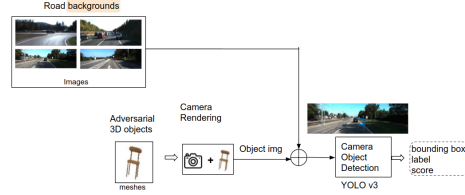
To address the challenges in Section 5.2, our test suite has the following designs;

Different driving backgrounds, 3D obstacle properties, and their interaction. To address C1, through our preliminary experiments, we choose driving backgrounds with different degrees of crowdedness, The 3D obstacle with various colors, shapes, and textures, and the relative position between the background and the obstacle. Fig. 5a is the overview of the automated test suite pipeline for the synthesized scenes.

For road backgrounds, 5 different driving scenarios are chosen according to the number of cars, the light condition, the shape of the road, and the emptiness of the road. As we can



(a) Overview of the automated test suite pipeline for benign synthesized scenes.



(b) Overview of the automated test suite for adversarial synthesized scenes.

Figure 5: Detection rate and confidence rate for benign objects

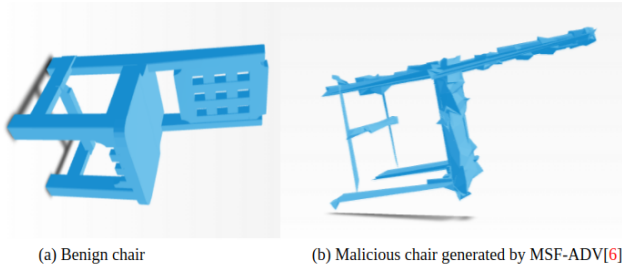


Figure 6: Benign chair and the malicious version

notice, some background has more than 13 cars while for others there are only 2 5 cars. The number of cars will influence the performance of the object detection neural network as it’s easier to have occlusion between the obstacle with the cars in the background when there are more cars. The result of the neural network will be quite sensitive to the position of the obstacle when the road is more crowded. These 5 backgrounds are with different light conditions also. For some backgrounds, the shadow of the trees takes up half part of the road while in other backgrounds there is less shadow on the road and even now shadows. The shadow area will affect the object detection accuracy in a way that obstacles put in the shadow are harder to be detected, especially when the color of the obstacle is quite dark.

For 3D obstacle properties, we first choose different shapes of objects of the same type. Considering that the chair is quite common and easy to obtain in real life, we choose the chair as our target benign obstacle and also use 3 different shapes of the chair for generality. We also change the color and the facing angle of the chair. The original color is wood brown, however, this is easy not to be detected when it’s under shadow or close to the soil. Therefore we also choose bright colors, i.e., blue and pink for the chair. Also, we rotate the chair for random angles as the angles in which the chair is put will also affect the detection performance. A chair with a full view will be detected easily while a chair with only a side view is harder to be detected and recognized as the chair.

For the interaction between the obstacle and the background, we adjust the position of the chair on both the x-axis

and y-axis. For a chair that is close to the camera, it’s easier to be detected while for a chair that is far from the camera, it’s the opposite case. We implement the distance between the chair and the camera by adjusting its position on the x-axis. The position of the chair will affect the occlusion. For example, if the chair is on the left side of the road and there happens to be a car here. The chair may block the car and prevent it to be detected. However, if the chair is on the right side, there won’t be occlusion then.

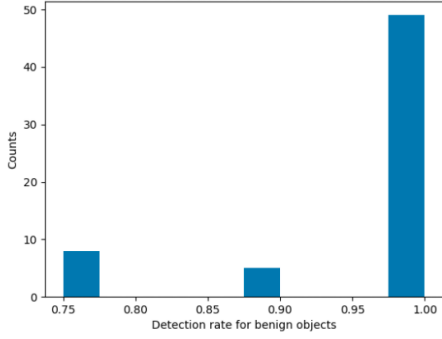
The attacker can easily simulate this by choosing an object with different shapes and colors and putting it on the road.

Adjust the size of the object along with its position. To address C2, we adopt camera imaging theory to adjust the height of the object so that it’s standing on the road. For example, if the obstacle is far from the camera, simply cropping the obstacle image and integrating it with a farther position on the x-axis won’t help. It’s likely to hit the road as the size of the obstacle should also shrink along with the increase of the distance to the camera. In the same way, if the obstacle is close to the camera, simply moving it closer on the x-axis will likely cause the object to float in the air. The size of the object should also increase when it’s closer to the camera. We do experiments on our backgrounds and take the following equation to model the relationship between the obstacle size and the distance between it and the camera.

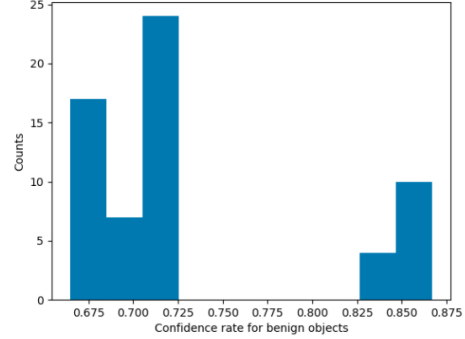
$$z = -1.73 + r/2 \quad (1)$$

Here z represents the distance between the obstacle and the camera in the unit of meters. r represents the scaling factor of the size of the obstacle. Moving the obstacle away or closer to the camera with the corresponding scaling factor in the equation ensures that the obstacle is standing on the ground without floating or hitting the ground.

On the other hand, we also consider the light condition to comply with the driving background. When there is a shadow in the background which is caused by the sunlight from a certain direction, an obstacle putting on the road should also have a similar shadow to ensure physical consistency. To produce a similar light condition, we consider the direction and intensity of the light when we’re using the camera rendering function to get the 2D image. Also, we start with normal ob-



(a) Detection rate for benign objects



(b) Confidence rate for benign objects

Figure 7: Detection rate and confidence rate for benign objects

stacles which can be obtained from life easily, e.g., a common chair. This is a fair assumption as getting the original obstacle shouldn't be difficult for the attacker.

Detection rate and Confidence rate To address **C3**, we parse the output of object detection neural networks, i.e., detected bounding box, object label, and confidence scores. Then we introduce two metrics, i.e., detection rate and confidence rate to measure the performance of the neural networks. Specifically, Fig. 10 shows the overview of the algorithm to get the total confidence scores and the correctly detected objects. Here we choose YOLOv3 as our target neural network as it's used in Autoware.AI [3] for object detection. The output of YOLOv3 includes the bounding box, label, and confidence score. The bounding box is used to describe the position of the detected object. The most important part of object detection is detecting the obstacle in the right place. Here we use the Intersection over Union(IoU) to describe the extent of overlap of two objects. Fig. 4 shows the equation to calculate IoU. We divide the area of overlap of two objects, i.e., the ground truth bounding box and the detected bounding box, with the area of the union of two objects. If the calculated IoU is close to 1, it means the detected bounding box almost overlaps with the ground truth and the detection position for the obstacle is correct. We set a threshold α to compare with the calculated IoU to decide if the detection position is correct.

Then the label of the obstacle is chosen and compared with the ground truth label. If they match, it means that the detected obstacle is correctly classified. And we think this obstacle is detected correctly. Therefore we record the confidence score of this detection result and count this as a correctly detected object.

In the end, we use the equation in Fig. 8 and Fig. 9 to calculate the detection rate and confidence rate. The detection rate is used to measure the percentage of correctly detected objects and the confidence rate is used to measure the average confidence scores per object. By comparing these two metrics in synthesized scenes with different settings, we can

$$\text{Detection rate} = \frac{\# \text{ of detected objects}}{\# \text{ of ground truth}}$$

Figure 8: Detection rate.

evaluate the performance of the neural network and therefore the authenticity of these scenes.

Automated pipeline for test suites To address **C4**, we need to design automated pipelines for the whole testing process to improve efficiency. In Fig. 5a, we first randomly choose one kind of chair as our target obstacle. Then we randomly choose the color between wood brown, blue and pink. The colored chair will be rotated at a random angle in $[0^\circ, 180^\circ]$. Meanwhile, a road background will be picked randomly to serve as the driving scenario. We also shift the position of the chair on both the x-axis and the y-axis. The shifting range in the x-axis is set to be in $[4\text{m}, 8.5\text{m}]$ and the shifting range in the y-axis is set to be in $[-8\text{m}, 8\text{m}]$. While we're shifting the obstacle, we also scale the size of the object according to the camera imaging equation to ensure physical consistency. Then we render this obstacle with a certain background to get the 2D images. This camera image is sent to the object detection neural network to get the detection rate and confidence rate of the final outputs. These are used to measure the authenticity of the benign synthesized image.

In Fig. 5b, we generate the adversarial obstacle using MSF-ADV [11]. Then we render the adversarial object with the camera rendering function and integrate it with the background. The generated 2D image is fed into the camera object detection neural network to get the detection rate and confidence rate of the adversarial synthesized scene. This is to measure the effectiveness of generated adversarial obstacles in the road background.

$$\text{Confidence rate} = \frac{\text{total confidence scores}}{\# \text{ of detected objects}}$$

Figure 9: Confidence rate.

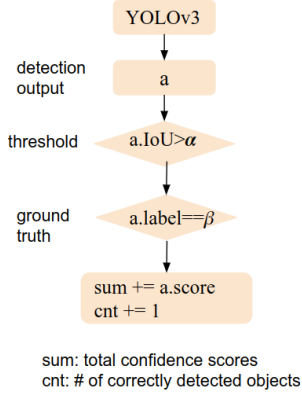


Figure 10: Algorithm for calculating confidence scores and detected objects.

5 Experiments

5.1 Evaluation Methodology and Setup

Object detection model selection In our evaluation, we target the object detection model provided in open-source industry-level AV systems to make our results more practical and realistic. Specifically, we choose the YOLOv3 [25] from the full-stack AV system, i.e., Autoware.AI [3] because Autoware.AI is applied in USDOT [5] and therefore the object detection neural network they’re using is representative. Besides, Autoware.AI has been installed in real physical driving vehicles and provided service on public roads [2]. This makes this test also practical in the real physical world. YOLOv3 also has a good performance as a real-time object detection algorithm and is widely used in the perception of AV system [3].

Synthesizing method selection Considering that there are many implementations for synthesizing the rendered obstacle with the road background, we experiment with the synthesizing pipeline used in MSF-ADV [11]. MSF-ADV [11] aims at designing adversarial obstacles which can fail both the camera and the LiDAR neural network detection. Since it takes thousands of optimization rounds to generate the adversarial objects, driving the physical vehicle on the road to get the camera and LiDAR sensor outputs with the updated obstacle in each optimization iteration is impractical. As a result, they design a way to synthesize the attack-influenced physical world digitally to get the resulted camera images and LiDAR point clouds. The adversarial obstacle generated with the digitally synthesized scene is also tested in the physical world and proved to be effective. This shows that their syn-

thesizing method might be realistic enough to simulate the physical world effect. Thus in our evaluation, we target the synthesizing method designed in MSF-ADV [11].

Adversarial attack selection Many prior works propose the physical adversarial attacks in AV systems, such as putting the drones flying in front of the target vehicle [35], placing an adversarial object on the rooftop of the target vehicle to hide this vehicle from the LiDAR detector [28]. However, due to the MSF design in current AV systems, which can recover the correct sensor reading as long as there is at least one sensor available, these adversarial attacks will fail in this setting. Thus, MSF-ADV [11] proposes a way to attack all fusion sensor sources at the same time. This is one of the most powerful attacks in perception modules in AV. Thus, in our evaluation, we will generate adversarial objects with the optimization-based method proposed by MSF-ADV [11].

3D object selection

Due to the practical object types for the camera models, we experiment with 3 different shapes of chairs from McGill 3D Shape Benchmark [7]: (1) a chair of size 0.6m * 0.5m * 1.5m, (2) a chair of size 0.4m * 0.7m * 1.3m, and (3) a chair of size 0.7m * 0.6m * 1.1m. These chairs are represented in the form of 3D point cloud mesh. Because the number of points is quite huge and brings a large overhead to the camera rendering process, we decrease the number of vertex and faces in the original 3D meshes with MeshLab [14]. It supports face reduction with high-quality preservation. We also scale the size of the chair according to the space in the background with MeshLab so it’s more realizable in real life.

Driving scenario selection For each chair with a certain shape, we select 5 real-world driving scenarios from the KITTI dataset [11]. Each driving scenario consists of the camera image, LiDAR point cloud, and the calibration matrix for each frame of sensor readings. KITTI offers various driving scenarios including different vehicles and roads with various properties.

5.2 Authenticity of Benign Synthesized Scenes

In this section, we evaluate the authenticity of the synthesized scene by measuring the performance of the object detection neural network.

Evaluation metrics. Given a randomly generated obstacle, driving background, and the relative position, we feed the rendered obstacle and the synthesized driving scene into the YOLOv3 model and test whether the cars in the original background and the newly added obstacle can be detected by YOLOv3. We use detection rate and confidence rate as metrics to measure the performance of YOLOv3 in the synthesized scene. Under this criterion, YOLOv3 should first successfully detect the obstacle in the blank background and the cars in the original driving background. In other words, if YOLOv3 detects the original obstacle and cars in the background separately in the first place, it means the obstacle and

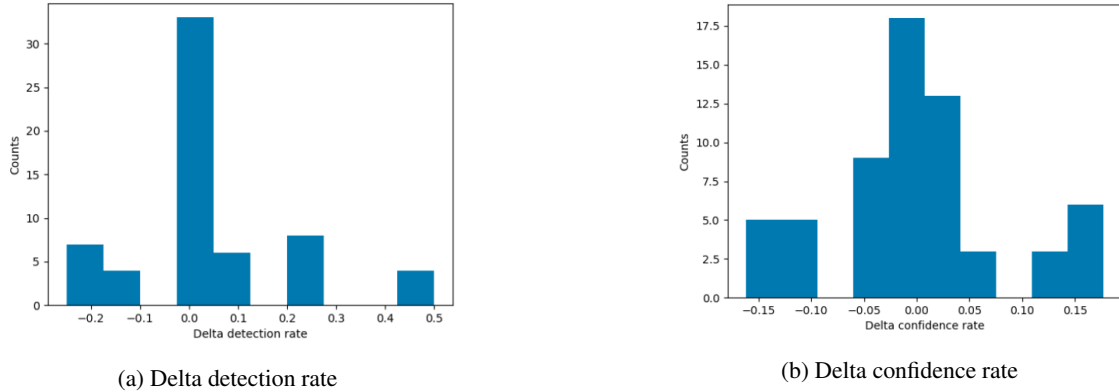


Figure 11: Delta detection rate and confidence rate between benign objects and adversarial ones.

cars are benign to YOLOv3. After synthesizing the obstacle with the driving background, a realistic scene should still allow YOLOv3 to detect the chair and equal amount of cars as it does in separate cases. A benign chair is the one without attack and thus should certainly be detected no matter what the background is.

Results.

Finding 1: The benign obstacle (i.e. chair) fails to be detected in all the tested settings. It means that when integrating the obstacle in the background, it's quite hard for the YOLOv3 model to detect it. This might be reasonable because perhaps the training dataset, didn't cover a similar image and the model hasn't seen a chair placed on the road. Therefore, if we want to simulate putting any obstacle on the road, we may have to choose the obstacle, background, and model carefully. As a result, the synthesized scene is not realistic enough especially for the chair, because the benign chair is never detected. *Finding 2: The benign obstacle may also lead to the performance decline in neural networks(e.g., occlusion).* Fig. 7a and Fig. 7b are the histograms for detection rate and confidence rate in benign objects. We draw this with more than 60 test settings. In Fig. 7a, for most cases, the detection rate is close to 1, which means it correctly detects all the cars in the background. However, there are several cases where YOLOv3 fails to detect all the cars. This is because the chair might block the car and as the result, the car is invisible to the neural network due to the occlusion.

5.3 Effectiveness of Adversarial Attacks on Malicious Scenes

Evaluation metrics. Given an obstacle, we generate the adversarial version with MSF-ADV [11] method, feed the rendered obstacle and the synthesized driving scene into the YOLOv3 model, and test whether the cars and the adversarial obstacle can be detected by YOLOv3. We use detection rate and confidence rate as metrics to measure the performance of YOLOv3 in the adversarial synthesized scene. Under this

criterion, YOLOv3 should first successfully detect the cars in the synthesized background but fail to detect the adversarial obstacle. This is to measure the effectiveness of adversarial attacks in the synthesized scene.

Results.

Finding 1: When the NN fails to detect the benign obstacle (i.e. chair), it's also hard to generate effective adversarial obstacles Fig. 11a and Fig. 11b are the histogram for delta detection rate and confidence rate. The delta value is calculated by subtracting the result in benign cases from the result in adversarial cases. If the delta value is close to 0, it means that the adversarial object doesn't have much influence on the detection accuracy. In Fig. 11a, we find for a large proportion of settings, the delta value is close to 0. We also find cases where the delta value is negative, which means that the adversarial object even improves the detection accuracy. In this case, there are more than 13 cars in the background chosen and they occlude each other a little bit. As the result, the detection accuracy might be very sensitive to the position and size of the chair. For the adversarial object, it may by accident improve the detection accuracy. The same goes with Fig. 11b, in which the negative delta value only appears when the background has countless cars.

6 Related Work

Adversarial camera and LiDAR-based attacks AVs Attacks in perception sensors can be divided into two categories, camera-based attack, and object-based attack. The camera-based attack methods [13, 17, 33] propose to hide the objects to be detected by adding adversarial patches. The attacker can apply different interference methods to enhance the robustness so that the objects won't be detected from varying observation angles and distances. This camera-based attack aims to change the texture of the object [11]. The Lidar-based attack methods [10, 12, 28, 35] propose to spoof the LiDAR with injecting laser [12], finding vulnerable LiDAR detection

locations [35] or changing the shape of the 3D objects [10]. This kind of attack can fool the LiDAR object detection mechanism, but it's hard to spoof cameras as it aims to change the shapes instead of the texture of the object [11]. In these works, to mislead the neural network, some outstanding patterns are generated to cause the model to have a tendency towards specific outputs.

Defense towards the adversarial camera and LiDAR-based attacks in AVs Defenses against these adversarial perception attacks also fall into two types. One kind of defense [29, 31, 34] aims to detect and recover the corrupted objects before they're sent to the detection algorithm. The authors reconstruct the objects with implicit functions [29] or denoising and upsampling [34]. Although these methods can achieve a good recovering rate, they focus on either camera-based attacks or object-based attack. The other kind of defense aims to fuse multiple sensors [20, 23, 24, 30] to avoid the spoofed sensor guiding the detection output. These Multiple Sensor Fusion (MSF) algorithms integrate the image and LiDAR feature map strategically to rely on the unattacked sensors.

References

- [1] The 6 levels of vehicle autonomy explained. <https://www.synopsys.com/automotive/autonomous-driving-levels.html>. Accessed: 2022-03-05.
- [2] Autoware self-driving vehicle on a highway. https://www.youtube.com/watch?v=npQMzH3j_d8. Accessed: 2022-03-05.
- [3] Autoware.ai. <https://www.autoware.org/>. Accessed: 2022-03-05.
- [4] Baidu apollo. <https://github.com/ApolloAuto/apollo>. Accessed: 2022-03-05.
- [5] carma-platform. <https://github.com/usdot-fhwa-stol/carma-platform>. Accessed: 2022-03-05.
- [6] Companies are racing to make self-driving cars. but why? <https://www.washingtonpost.com/outlook/2022/02/04/self-driving-cars-why>. Accessed: 2022-03-05.
- [7] McGill 3d shape benchmark. <http://www.cim.mcgill.ca/~shape/benchMark/>. Accessed: 2022-03-05.
- [8] Ray tracing. <https://computationalthinking.mit.edu/Fall20/lecture14/>. Accessed: 2022-03-05.
- [9] Taxonomy and definitions for terms related to driving automation systems for on-road motor vehicles.
- [10] Mazen Abdelfattah, Kaiwen Yuan, Z. Jane Wang, and Rabab Kreidieh Ward. Towards universal physical attacks on cascaded camera-lidar 3d object detection models. In *ICIP*, 2021.
- [11] Yulong Cao, Ningfei Wang, Chaowei Xiao, Dawei Yang, Jin Fang, Ruigang Yang, Qi Alfred Chen, Mingyan Liu, and Bo Li. Invisible for both camera and lidar: Security of multi-sensor fusion based perception in autonomous driving under physical-world attacks. *Proceedings - IEEE Symposium on Security and Privacy*, May 2021.
- [12] Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Wonseok Park, Sara Rampazzi, Qi Alfred Chen, Kevin Fu, and Z. Morley Mao. Adversarial sensor attack on lidar-based perception in autonomous driving. *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019.
- [13] Shang-Tse Chen and Jason Martin. Physical adversarial attack on object detectors (extended abstract). 2018.
- [14] Paolo Cignoni, Marco Callieri, Massimiliano Corsini, Matteo Dellepiane, Fabio Ganovelli, and Guido Ranzuglia. MeshLab: an Open-Source Mesh Processing Tool. In Vittorio Scarano, Rosario De Chiara, and Ugo Erra, editors, *Eurographics Italian Chapter Conference*. The Eurographics Association, 2008.
- [15] Ronald T. H. Collis. Lidar. *Encyclopedic Dictionary of Archaeology*, 2021.
- [16] Object Detectors and In The. Camou: Learning a vehicle camouflage. 2018.
- [17] Kevin Eykholt, I. Evtimov, Earlece Fernandes, Bo Li, Amir Rahmati, Florian Tramèr, Atul Prakash, Tadayoshi Kohno, and Dawn Xiaodong Song. Physical adversarial examples for object detectors. *ArXiv*, abs/1807.07769, 2018.
- [18] Kevin Eykholt, I. Evtimov, Earlece Fernandes, Bo Li, Amir Rahmati, Florian Tramèr, Atul Prakash, Tadayoshi Kohno, and Dawn Xiaodong Song. Physical adversarial examples for object detectors. *ArXiv*, abs/1807.07769, 2018.
- [19] Davi Frossard and Raquel Urtasun. End-to-end learning of multi-sensor 3d tracking by detection. *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 635–642, 2018.
- [20] Davi Frossard and Raquel Urtasun. End-to-end learning of multi-sensor 3d tracking by detection. *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 635–642, 2018.

- [21] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013.
- [22] Hiroharu Kato, Y. Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3907–3916, 2018.
- [23] Ming Liang, Binh Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7337–7345, 2019.
- [24] Ming Liang, Binh Yang, Shenlong Wang, and Raquel Urtasun. Deep continuous fusion for multi-sensor 3d object detection. In *ECCV*, 2018.
- [25] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *ArXiv*, abs/1804.02767, 2018.
- [26] Jiachen Sun, Yulong Cao, Qi Alfred Chen, and Z. Morley Mao. Towards robust lidar-based perception in autonomous driving: General black-box adversarial sensor attack and countermeasures. *ArXiv*, abs/2006.16974, 2020.
- [27] Jiachen Sun, Yulong Cao, Qi Alfred Chen, and Z. Morley Mao. Towards robust lidar-based perception in autonomous driving: General black-box adversarial sensor attack and countermeasures. *ArXiv*, abs/2006.16974, 2020.
- [28] James Tu, Mengye Ren, Sivabalan Manivasagam, Ming Liang, Binh Yang, Richard Du, Frank Cheng, and Raquel Urtasun. Physically realizable adversarial examples for lidar object detection. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13713–13722, 2020.
- [29] Ziyi Wu, Yueqi Duan, He Wang, Qingnan Fan, and Leonidas J. Guibas. If-defense: 3d adversarial point cloud defense via implicit function based restoration. *ArXiv*, abs/2010.05272, 2020.
- [30] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Point-fusion: Deep sensor fusion for 3d bounding box estimation. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2018.
- [31] Jiancheng Yang, Qiang Zhang, Rongyao Fang, Bingbing Ni, Jinxian Liu, and Qi Tian. Adversarial attack and defense on point sets. *ArXiv*, abs/1902.10899, 2019.
- [32] Yue Zhao, Hong Zhu, Ruigang Liang, Qintao Shen, Shengzhi Zhang, and Kai Chen. Seeing isn’t believing: Practical adversarial attack against object detectors. *arXiv: Computer Vision and Pattern Recognition*, 2018.
- [33] Yue Zhao, Hong Zhu, Ruigang Liang, Qintao Shen, Shengzhi Zhang, and Kai Chen. Seeing isn’t believing: Towards more robust adversarial attack against real world object detectors. *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019.
- [34] Hang Zhou, Kejiang Chen, Weiming Zhang, Han Fang, Wenbo Zhou, and Nenghai Yu. Dup-net: Denoiser and upsampler network for 3d adversarial point clouds defense. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1961–1970, 2019.
- [35] Yi Zhu, Chenglin Miao, T. Zheng, Foad Hajiaghajani, Lu Su, and Chunming Qiao. Can we use arbitrary objects to attack lidar perception in autonomous driving? *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021.