
Project Proposal - Attacks on Automated Vehicles

Connor Raymond Stewart

University of Waterloo

200 University Ave W, Waterloo, ON N2L 3G1

crstewar@uwaterloo.ca

Abstract

Many recent advances in machine learning are being used to develop AI-piloted, driver-less autonomous vehicles. Despite recent advances in machine learning and AI approaches to autonomous vehicle development, many issues related to autonomous vehicles remain, such as the problems related to the robustness of machine learning models. Planning and perception modules in autonomous vehicles can be fooled into making incorrect inferences about the environment given slight modifications to vehicular sensory input. Therefore, the project aims to research a set of methods to test and evaluate model robustness to conduct a theoretical evaluation to determine the best methods to develop robust autonomous vehicle models.

1 Background

Multiple advances in machine learning have been made, which has helped expand autonomous self-driving AI vehicles' ability to reach a level of autonomy where human intervention on the road is no longer needed. However, safety concerns, including adversarial attacks on both perception-based and planning-based models, have been observed. The following report introduces how machine learning modules are fooled to create a theoretical evaluation of such models.

1.1 Perception

Perception models are susceptible to attack through slight perturbations in image data, and small perturbation on image data has been shown to cause erroneous output from a trained DNN model [1]. Physical attacks like placing black lines on stop signs and roads have been shown to cause erroneous output - which leads to actual movements - in self-driving autonomous vehicles [1]. Bayesian Optimisation with a robust objective function can be used to hijack an autonomous vehicle. It can cause it to deviate from its path onto a path crafted by an attacker [1]. Methods to combat these problems can be explored by analyzing errors associated with the networks, and testing data can be analyzed to determine correct defences.

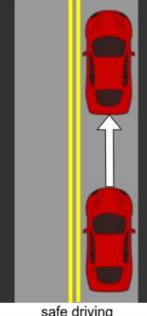
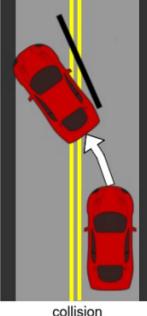
Domain	Input	Adversary	Result
Image		pixel level	
Physical		physical attack	
End-to-end Autonomous Driving	 safe driving	 physical attack	 collision

Figure 1: Vision-Based Perception can lead to physical attacks such as placing lines on objects - as seen in (a) - to fool vehicle perception modules into making inaccurate environmental classifications - as seen in (b) - which lead to collisions [1].

26 Adversarial attacks on sensor cameras and light detection and ranging systems are vulnerable to
 27 optimization methods for input perturbation [8]. Discussing defences against algorithmic methods
 28 and attempting to evaluate defences against formulas to fool vehicle sensors will be covered in the
 29 report.

30 Environment sensors such as dynamic vehicle sensors (Tire Pressure Monitoring Systems (TPMS),
 31 magnetic encoders, and inertial sensors) and environment sensors (Light Detection and Ranging
 32 (LiDAR), ultrasonic, camera, Radio Detection and Ranging (Radar) systems, and Global Positioning
 33 System (GPS) units) are attainable by cyber-attackers [9]. Attacks and defences for attacks are
 34 discussed in detail in the article published by *El-Rewin et al.*, and attack methods like spoofing and
 35 jamming are discussed along with machine-learning countermeasures [9].

36 Various solutions to prevent physical attacks on perception-based vehicle systems have been elaborated
 37 and include the development of active and passive algorithms to detect attacks and mitigate
 38 their effects on a system [6]. Exploring these topics is undertaken in the project to determine models
 39 that can be used to create machine learning models to control attacks [6].

40 Intelligently selected adversarial examples have been explored by *Chowdhury et al.* and shown to
 41 increase vehicle robustness [2]. Machine learning methods can monitor for abnormalities in-vehicle
 42 monitors, and robustness methods can defend against input irregularities [2].

43 Reasons behind perturbation attack effectiveness (ghost object detection) and mitigation strategies to
 44 prevent such attacks - such as using Spatial Pyramid Pooling - are elaborated upon [5].

45 1.2 Planning

46 PlanFuzz is a dynamic testing system used to discover vulnerabilities associated with vehicle planning
 47 module-based behaviours when presented with environmental obstacles on a roadway, such as bikes
 48 or boxes [10]. Adversarial examples such as roadside debris can cause a vehicle to accidentally
 49 detect a roadblock or road sign on the reads [10]. A system of constraints for the attacker-introduced
 50 physical objects (the planning invariant) is introduced [10]. Furthermore, a planning invariant aware
 51 physical object generation system to enforce driving rules and a behavioural planning vulnerability

52 distance system to measure how close the current planning decision is to violating PI and thus trigger
53 a vulnerability in the run time is introduced [10]. Considering that the vehicle's planning system
54 can be attacked by placing otherwise benign objects on the road in strategic locations, developing a
55 robust learning system which omits benign objects or plans around their existence is necessary to
56 prevent adversarial attacks on the planning module [10].

57 Fine-tuning the victim against the adversary by training against adversaries that engage in seemingly
58 random or bizarre behaviour can benefit self-driving cars. Adversarial examples on deep reinforce-
59 ment learning models show that zero-sum games can be lost by confusing an AI agent by having an
60 adversary play seemingly random choices [3]. If extended to autonomous vehicles, it is essential
61 to realize the problems with a driver on the road who intentionally behaves irrationally. A vehicle
62 cannot make irrational driving decisions in the face of an adversarial driver. The system presented in
63 the paper lays the groundwork for robustness against poor drivers on the road.

64 AVMaestro presents an efficient and effective policy framework to allow autonomous vehicles' safe
65 execution [11]. Systems to allow for detecting anomalies on the roads, along with defensive methods
66 to avoid bad decision-making in their presence, are discussed in detail [11].

67 1.3 Testing and Evaluation

68 The appropriate methods to plan, test, and debug autonomous vehicles with machine learning concepts
69 are characterized by a lack of agreement in the industry. A framework for testing and evaluating
70 machine learning systems is presented by *Tuncali et al.* by converting arrays to test combinations of
71 discrete parameters and simulating annealing to find corner cases [4]. Understanding more methods
72 to test robustness is vital to evaluating the performance of machine learning models to determine the
73 best methods to make vehicles more robust.

74 General defences against adversarial examples for deep neural networks and algorithms to defend
75 against adversarial attacks can be read in more detail in the paper written by *Papernot et al.* [7].

76 2 Introduction

77 In the project, theoretical analysis and or an algorithmic design are proposed. An theoretical analysis
78 inspired by and using the data from the proposed literature sources, along with an algorithmic
79 analysis consisting of a explanation of the various methods, will be written up and tested for the paper
80 summary and presentation. The project aims to study then which methods are most effective and
81 provide quantitative and qualitative analysis to determine how models operate. Results from papers
82 will be evaluated, and theoretical and algorithmic explanations of the results from authors will be
83 created.

84 Therefore, the project's purpose is not to produce a working adversarial defence system but to evaluate
85 existing defence systems to determine how well they perform in general. Code sources from GitHub
86 produced by article authors will be evaluated and reproduced, and tests will be run to determine
87 approaches to machine learning robustness evaluation.

88 Results will include explanations to empirical data generated from literature sources to determine
89 how machine learning robustness methods compare to one another in the autonomous vehicle system
90 context. Algorithms are elaborated on past how they are explained in articles, and an algorithmic
91 analysis may be undertaken to elaborate on the best process discovered for evaluating or enforcing
92 robustness.

93 In the project herein, three problems will be evaluated. The first is the reproduction of and theoretical
94 analysis of previously created empirical results of perception-based attacks. The second problem is
95 an algorithmic analysis and explanation of methods for fooling self-driving vehicle planning modules
96 and testing their robustness rate. Finally, the third problem is a theoretical analysis of the different
97 systems for testing adversarial robustness.

98 The next section of the paper focuses on detailed information regarding the previous works used
99 for the project, with the following section focusing on the methodology behind the project's three
100 problems. Results of the evaluation performed from the empirical analysis will be presented in a final
101 section after the methodology.

102 **3 Previous Works**

103 Multiple topics are covered in the project, including perceptron module attacks, planning module
 104 attacks, and testing and evaluation of attacks. Perceptron module attacks focus exclusively on the
 105 perception, classification, and identification of objects by autonomous vehicle sensors. Topics include
 106 vision-based perceptron attacks, perturbation attacks, physical sensor attacks, and supporting topics.
 107 Planning module attacks focus on planning module vulnerabilities which can lead to abnormal
 108 driving decisions based on adversarial vehicle behaviour, planning around physical obstacles in
 109 an autonomous vehicle driving environment, and adversarial policy attacks. Finally, testing and
 110 evaluation of attacks include the methodology behind the planning, testing, and debugging of attacks
 111 on autonomous vehicles and general defences for neural networks against adversarial inputs.

112 **3.1 Perceptron Module Attacks**

113 The perceptron module can be attacked to fool a vehicle into misinterpreting sensory input. Such
 114 misinterpretations can cause vehicles to make situationally inappropriate choices. We will investigate
 115 algorithms and data evaluating such attacks.

116 **Vision-Based Perceptron Attack** Several objective functions are used to evaluate the algorithms
 117 presented in the vision-based perception system presented by *Boloor et al.* Essentially, attacks are a
 118 set of length, width, and rotation angles, and attacks are generated through random or grid search
 119 functions which are evaluated with objective functions.

Algorithm 1 Adversary Search Algorithm

```

Require: Strategy ∈ Random, Grid
     $i \leftarrow 0$ 
    MetricsList  $\leftarrow [ ]$ 
    loop
         $\delta_i \leftarrow \text{GenerateAttack}(\text{Strategy})$ 
        results  $\leftarrow \text{RunScenario}(\delta_i)$ 
         $y_i \leftarrow \text{CalculateObjectiveFunction(results)}$ 
        MetricsList.append( $y_i$ )
         $i \leftarrow i + 1$ 
    end loop
    return arg max MetricsList

```

Figure 2: The RunScenario function runs a simulation and returns vehicle data, and the CalculateObjectiveFunction calculates an objective function with the results of the RunScenario function [1]

120 Larger pattern spaces benefit enormously from Bayesian optimization since it does not require
 121 gradient information.

$$z = \frac{\mu_{f|\mathcal{D}}(\delta) - y_{\max}}{\sigma_{f|\mathcal{D}}(\delta)};$$

$$u(\delta) = (\mu_{f|\mathcal{D}}(\delta) - y_{\max})\Phi(z) + \sigma_{f|\mathcal{D}}(\delta)\phi(z),$$

Figure 3: The authors presented the following closed-form solution. [1]

122 Bayesian optimization's lower runtime complexity leads the authors to create a second algorithmic
 123 method for determining an adversary search.

Algorithm 2 Bayesian Adversary Search Algorithm

```

i  $\leftarrow 0$ 
MetricsList  $\leftarrow [ ]$ 
loop
 $\delta_i \leftarrow \arg \max u(\delta)$ 
results  $\leftarrow \text{RunScenario}(\delta_i)$ 
 $y_i \leftarrow \text{CalculateObjectiveFunction}(\text{results})$ 
MetricsList.append( $y_i$ )
Update Gaussian Process and  $\mathcal{D}$  with  $(\delta_i, y_i)$ 
i  $\leftarrow i + 1$ 
end loop
return  $\arg \max \text{MetricsList}$ 

```

Figure 4: The optimized adversarial search algorithm [1]

124 The authors use objective functions based on several parameters. The objective function used in the
 125 algorithmic implementations by the authors are presented below:

$$\vec{\Theta}_\delta = [\theta_{F_l}, \theta_{F_{l+1}}, \dots, \theta_{F_{l+\Delta}}]$$

Figure 5: Predicted Steering Vector [1]

$$\vec{p}_\delta = [(x_l, y_l), (x_{l+1}, y_{l+1}), \dots, (x_{l+\Delta}, y_{l+\Delta})]$$

Figure 6: Change in position (x, y) across frames in a given attack [1]

$$\begin{aligned} \text{Collide Right : } & \max_{l, \delta} \sum_{i=0}^{\Delta} \vec{\Theta}_{\delta_i} \\ \text{Collide Left : } & \min_{l, \delta} \sum_{i=0}^{\Delta} \vec{\Theta}_{\delta_i} \\ \text{subject to : } & l \in L, \quad \delta \in S. \end{aligned}$$

Figure 7: Steering angle summation objectives for the min-max sum of steering angles [1]

$$\begin{aligned} & \max_{l, \delta} \|\vec{\Theta}_\delta - \vec{\Theta}_{\text{baseline}}\|_1 \\ \text{subject to : } & l \in L, \quad \delta \in S. \end{aligned}$$

Figure 8: Absolute Steering Angle Differences for the most significant change in steering angles [1]

$$\begin{aligned} & \max_{l, \delta} \|\vec{p}_\delta - \vec{p}_{\text{baseline}}\|_2 \\ \text{subject to : } & l \in L, \quad \delta \in S. \end{aligned}$$

Figure 9: Path Deviation for optimizing deviation from baseline path making [1]

126 We make a note of the following empirical data as presented by the authors. When modifying the
 127 sum of steering angles, we notice that the safe driving rate only reduces significantly for left turns.

128 Path deviation works best at reducing safe driving for straight car movements, and only halves safe
 129 driving for left and right turns. The absolute steering difference parameter almost removes all safe
 130 driving for left and right movements and strongly reduces safe driving for right movements. The
 131 absolute standard steering difference generates a very high collision rate for right-hand moves.

Metric	Left			Straight			Right					
	safe	collision	offroad	opp. lane	safe	collision	offroad	opp. lane	safe	collision	offroad	opp. lane
\sum st. angles	18.2	0	0	81.8	99	0	0	1	72.2	9.5	13.8	24.5
path deviation	64.6	0	0	35.4	23.8	2.5	2.8	76.2	57.2	24.0	28.3	40.2
abs. st. diff.	0.2	0	0	99.8	22.7	7.5	9.3	77.3	0	95.2	99.2	100

Figure 10: Comparing candidate objective functions shows that absolute standard difference works the best [1]

132 Results show that Bayesian optimization finds a more significant number of unique successful
 133 adversaries in the same number of iterations as random and naive grid searches but also finds these
 134 adversarial attacks faster.

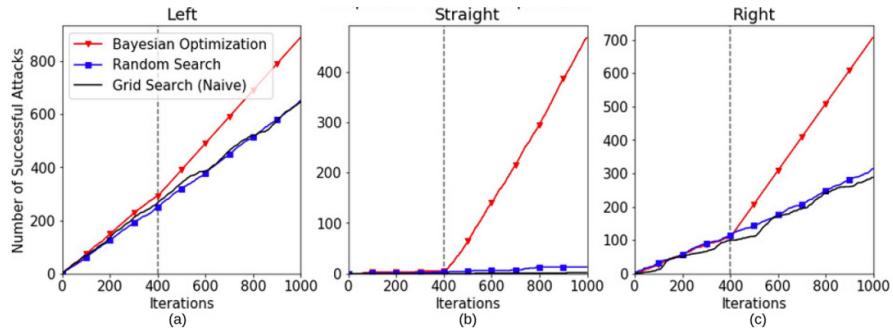


Figure 11: Bayesian optimization (red) shows strong results when compared to other search optimizers [1]

135 Finally, when measuring adversaries against right-turn driving, the number of infractions is propor-
 136 tional to the change in the summation of steering angles. Furthermore, we notice patterns that cause
 137 the minimum steering sum and maximum collision to look similar. Therefore, there is a pattern
 138 between minimizing the steering sum and causing the autonomous vehicle to make an infraction.

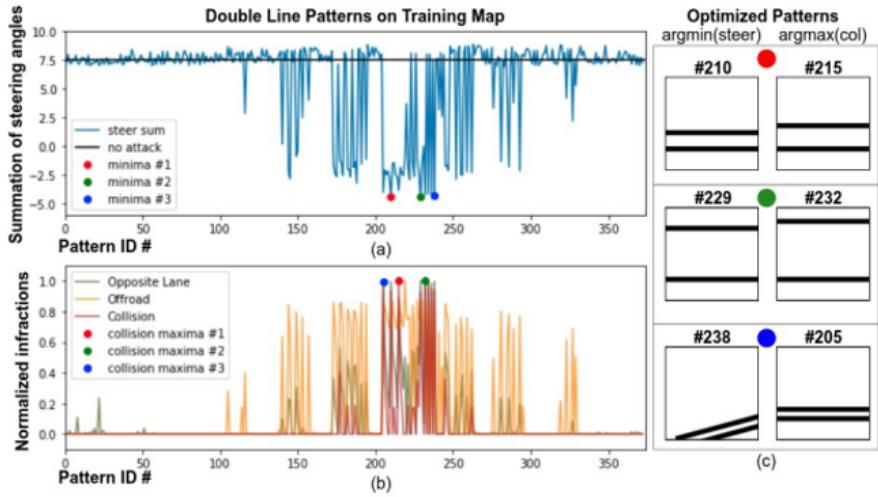


Figure 12: Adversary analysis against right-turn driving to compare summation of steering angles with the normalized infractions shows that patterns which cause minimum steering sum and maximum collisions are similar [1]

139 **Perturbation Attacks** The algorithm presented by Zhang *et al.* attempts to add pixel sized changes
 140 to make deep learning models dysfunctional is the purpose of the system [5]. We can consider a
 141 system where a perturbation is updated using projected loss gradient of a 3D object detector via
 142 multiple iterations.

$$\begin{aligned}\delta_n^{\text{per}} &= \text{Clip}_{\epsilon} \{ \alpha \times \text{sign}(\nabla_{(I_l, I_r)} L(O_\theta(\tilde{I}_{l,n}^{\text{per}}, \tilde{I}_{r,n}^{\text{per}}), b^{\text{true}})) \} \\ (\tilde{I}_{l,n+1}^{\text{per}}, \tilde{I}_{r,n+1}^{\text{per}}) &= (\tilde{I}_{l,n}^{\text{per}} + \delta_n^{\text{per}}, \tilde{I}_{r,n}^{\text{per}} + \delta_n^{\text{per}})\end{aligned}$$

Figure 13: Projected loss gradient of a 3D object detector via multiple iterations [5]

143 Furthermore, we can represent a patched image pair as follows:

$$\underset{\delta^{\text{pat}}}{\text{argmin}} \mathbb{E}_{(I_l, I_r) \sim \mathcal{I}, (loc_l, loc_r) \sim \mathcal{L}, \tau \sim \mathcal{T}} L(O_\theta(\tilde{I}_l^{\text{pat}}, \tilde{I}_r^{\text{pat}}), b^*)$$

Figure 14: Patched image pair minimization [5]

144 As can be seen in the below table, the driving safety performance metric under the perturbation attack
 145 reveals that the average success rate is lower as the iteration value and alpha value increase. Stereo
 146 R-CNN models have a noticeably more significant change of a collision on the right turns when
 147 attacked when compared with the DSGN model, when alpha equals zero-point-four. We also notice
 148 that the Stereo R-CNN model has a noticeably greater chance of collision on the right, left, and
 149 straight vehicle movements when attacked compared with the DSGN model, when alpha equals one.

150 We note that alpha represents the attack intensity, n represents the number of iterations, the iteration
 151 counter represents the planned trajectory, DGSN means Deep Stereo Geometry Network, and R-CNN
 152 means Region-Based Convolutional Neural Network [5]. Furthermore, we notice that the success
 153 rate represents the percent chance of success that a vehicle trajectory is generated, the collision rate
 154 represents the percent chance that a trajectory leads to a collision, and the safe driving rate is the
 155 percentage chance that a trajectory leads to a collision-free outcome [5].

DRIVING SAFETY PERFORMANCE METRICS UNDER THE PERTURBATION ATTACK ($\alpha = 1$)													
Model	Iteration	DSGN				Stereo R-CNN				Unattacked	2	3	4
		Left	Unattacked	1	2	3	4	Unattacked	1				
Success rate (%)	Straight	89.6	90.0	89.3	89.8	89.5	89.8	90.2	88.4	89.8	84.6	95.8	94.2
	Left	96.5	96.7	96.5	96.5	96.5	96.9	96.4	95.8	94.2	95.5	90.6	90.6
	Right	84.4	84.4	84.4	84.6	84.8	84.9	85.7	85.1	86.7	82.4	82.4	82.4
Collision rate (%)	Left	2.2	2.2	2.3	2.4	2.4	2.4	3.1	3.2	3.2	3.4	3.4	3.4
	Straight	0.7	0.7	0.9	1.3	1.3	1.1	1.5	1.2	2.1	1.5	1.5	1.5
	Right	1.7	3.2	3.8	5.3	5.9	3.3	4.1	5.1	8.6	11.3	11.3	11.3
Safe driving rate (%)	Left	87.7	87.7	86.7	86.7	86.4	87.0	87.2	85.5	85.5	79.2	87.4	87.4
	Straight	95.8	96.0	95.6	92.9	95.6	95.3	93.4	93.3	93.4	87.4	93.4	93.4
	Right	83.0	81.7	81.1	80.1	79.7	82.1	82.1	80.7	79.2	73.0	87.4	87.4

DRIVING SAFETY PERFORMANCE METRICS UNDER THE PERTURBATION ATTACK ($\alpha = 0.4$)													
Model	Iteration	DSGN				Stereo R-CNN				Unattacked	2	3	4
		Left	Unattacked	1	2	3	4	Unattacked	1				
Success rate (%)	Straight	89.6	89.6	89.5	89.3	89.5	89.8	89.8	90.4	90.0	90.8	95.0	94.5
	Left	96.5	96.7	96.7	96.7	96.7	96.9	96.4	96.0	96.1	95.1	94.5	94.5
	Right	84.4	84.2	84.4	84.6	84.8	84.9	85.3	84.4	85.9	85.2	85.2	85.2
Collision rate (%)	Left	2.2	2.2	2.4	2.6	2.8	3.1	2.7	3.7	4.1	3.7	3.7	3.7
	Straight	0.7	0.7	0.9	1.1	1.1	1.1	1.1	0.9	1.3	1.7	1.7	1.7
	Right	1.7	2.7	3.6	4.8	5.8	3.3	3.5	4.4	4.2	4.2	4.2	4.2
Safe driving rate (%)	Left	87.7	87.7	87.3	86.9	86.9	87.0	87.4	87.0	86.3	85.4	85.4	85.4
	Straight	95.8	96.0	96.0	95.8	95.8	95.3	94.9	94.1	93.8	92.9	92.9	92.9
	Right	83.0	82.2	82.0	81.5	81.5	82.1	82.3	81.4	81.2	79.1	79.1	79.1

Figure 15: Driving Safety Rates when under Perturbation Attacks [5]

156 In the tables below, it is essential to note that AP is an abbreviation for average precision [5]. The
157 scenario represents the difficulty of a benchmark, random attacks are performed by placing trained
158 patches at a random position in the image, and specific attacks are performed by placing trained
159 patches randomly within a specific region of the image [5].

AVERAGE PRECISION FOR 3D OBJECT DETECTION UNDER PATCH ATTACK														
Model	Scenario	Iteration	DSGN				Stereo R-CNN				Unattacked	2	3	4
			Random	Attack	Left	Specific	Random	Attack	Left	Specific				
AP (%)	3D	easy	70.94	63.85	63.85	60.9	60.59	64.87	56.47	53.17	48.14	73.82	70.07	70.07
	detection	moderate	52.98	48.20	51.47	50.77	51.63	38.20	37.07	36.27	35.23	38.02		
	hard	47.29	44.30	46.80	46.15	46.48	32.66	31.88	31.21	30.60	32.45			

DRIVING SAFETY PERFORMANCE METRICS UNDER THE PATCH ATTACK														
Model	Scenario	Iteration	DSGN				Stereo R-CNN				Unattacked	2	3	4
			Random	Attack	Left	Specific	Random	Attack	Left	Specific				
Success rate(%)	Straight	89.6	90.1	90.1	90.6	90.6	89.8	89.8	90.3	90.3	83.7			
	Left	96.5	96.5	96.5	96.7	96.7	96.4	96.4	91.2	91.2	57.5			
	Right	84.4	84.6	84.6	84.8	84.8	94.9	94.9	68.7	74.4				
Collision rate(%)	Left	2.2	2.2	2.8	2.8	2.8	3.1	3.1	1.9	1.9	4.1			
	Straight	0.7	0.7	0.9	1.1	1.1	1.1	1.1	1.4	4.7				
	Right	1.7	2.3	2.3	2.1	2.1	3.3	3.3	3.9	4.2				
Safe driving rate(%)	Left	87.7	87.7	87.5	87.5	87.5	87.0	87.0	88.8	88.8	81.0			
	Straight	95.8	95.8	95.8	95.8	95.8	95.3	95.3	80.1	84.7				
	Right	83.0	82.6	83.0	82.1	82.1	65.9	65.9	71.2					

Figure 16: Driving Safety Performance Metrics under Patch Attacks [5]

160 We notice that random and specific attacks affect the average precision for 3D object detection under
161 patch attacks. Hard scenarios appear to be much more robust than easy and moderate scenarios.
162 DSGN models appear to be more robust to specific attacks, whereas stereo R-CNN models are more
163 robust to random attacks. Random and specific attacks do not affect the safe driving rate of DSGN
164 models very much, but random and specific attacks appear to increase the collision rate for left and
165 right turns when used against DSGN models. Notably, Stereo R-CNN models show a dramatic
166 reduction in safe driving when turning left or right under random attacks or turning right during
167 specific attacks. We also notice a significant drop in safe driving rate under both random and specific
168 attacks for straight vehicle movements. A moderate reduction in safe driving when left turning under
169 specific attacks is also noted. The Stereo R-CNN model shows a significant increase in collision rate
170 under specific attacks when moving straight or right. Interestingly, under both attacks, the R-CNN
171 model's collision rate drops. The DSGN model is comparatively more resistant to random attacks
172 than specific attacks, whereas the Stereo R-CNN model is comparatively more resistant to specific
173 attacks when compared to random attacks.

174 In the top table below, specific attack scenarios are very close to the same intention cases, whereas
175 different intentions are close to attacked scenarios [5]. In the middle diagram, ground truth indicates
176 if an object is moving [5]. A patch attack significantly affects a small number of pixels, whereas a
177 perturbation attack slightly affects many pixels [5]. In the mid-table, DSGN models are very robust
178 under perturbation and patch attacks, showing only slight drops in safe driving rates under all attacks.
179 The Stereo R-CNN model shows large drops in safe driving rates under perturbation attacks and
180 substantial drops in safe driving rates under patch attacks. In the final table, we notice that greedy
181 breadth-first searches (GBFS) have lower ground truth safe driving rates than A* planning algorithms
182 when planning left, and right turns. Ground truth data of 3D object detection is used as inputs
183 for motion planning algorithms [5]. Therefore, the planning algorithm can affect the self-driving
184 vehicle's safe driving rate and robustness.

DRIVING SAFETY PERFORMANCE METRICS OF STEREO R-CNN UNDER THE PATCH ATTACK WITH VARIOUS INTENTIONS

	Specific attack		Random attack		
	Same intentions ¹	Different intentions ²	Same intentions ¹	Different intentions ²	Unattacked
Success rate (%)	71.0	76.2	91.6	90.4	
Collision rate (%)	3.5	2.3	1.5	2.4	
Safe driving rate (%)	69.3	74.4	90.1	88.1	

¹ "Same intentions" refers to cases where the attack intention and the driving intention are the same.² "Different intentions" refers to cases where the attack intention differs from the driving intention.

SAFE DRIVING RATE WITH DIFFERENT INPUTS

Model	DSGN			Stereo R-CNN			
	Ground Truth	Unattacked	Perturbation Attack	Patch Attack	Unattacked	Perturbation Attack	Patch Attack
Safe driving rate (%)	Left	89.0	87.7	86.7	87.0	79.2	68.9
	Straight	98.0	95.8	95.6	95.8	95.3	87.7
	Right	85.2	83.0	79.7	82.6	82.1	73.0
							65.9

SAFE DRIVING RATE USING DIFFERENT PLANNING ALGORITHMS

Planning algorithm	GIBFS	Δ^+	
Safe driving rate (%)	Left	87.9	89.7
	Straight	98.0	98.0
	Right	82.3	85.2

Figure 17: Safe Driving Rates under Various Types of Attacks [5]

185 **Physical Sensor Attacks** Understanding the types of defences against algorithmic methods and
 186 attempting to evaluate defences against formulas to fool vehicle sensors is required to prevent sensory-
 187 based attacks. We note that Cao *et al.* and Raj Gautam Dutta (2018) use complex derivations and
 188 linear algebraic models for camera input to model attacks in their algorithmic systems. Complex
 189 algorithms for attack generation and defence are created to prevent attacks on vehicular sensors.

Algorithm 1: Generating adversarial examples by leveraging global spatial transformation

```

input: Target model:  $M$ ;  

       3D point cloud  $X$  ;  

       3D spoofed 3D point cloud  $T$ ;  

 $\theta$  Optimizer  $opt$ ;  

       Max iteration  $N$ ;  

output: 3D adversarial 3D point cloud  $X'$ ;  

Initialization:  $\theta \leftarrow 0, \tau_x \leftarrow 0, s_h \leftarrow 1, l_{min} = +\infty,$   

 $x = \Phi(X), t = \Phi(T);$   

/* Initiate parameters by sampling around the  

   transformation parameters  $Target_\theta, Target_{\tau_x}$  that  

   transforms  $t$  to the target position  $(px, py)$  of  

   the attack */  

for  $i\tau_x \leftarrow -L_t$  to  $L_t$  do  

for  $i\theta \leftarrow -L_\theta$  to  $L_\theta$  do  

  /* Initialize parameter . */  

 $\theta \leftarrow Target_\theta + i\theta, \tau_x \leftarrow Target_{\tau_x} + \tau_x i;$   

for  $iter \leftarrow 1$  to  $N$  do  

  /* Calculate adversarial loss */  

 $l_{adv} \leftarrow$  Equation 7;  

  /* Update the parameters  $\theta, \tau_x, s_h$  based on  

   optimizer  $opt$  and loss  $l_{adv}$  */  

 $\theta, \tau_x, s_h \leftarrow opt(l_{adv}; \theta, \tau_x, s_h)$   

if  $l_{min} < l_{adv}$  then  

  |  $\theta^{final}, \tau_x^{final}, s_h^{final} \leftarrow \theta, \tau_x, s_h$   

end  

end  

end  

end  

 $T' \leftarrow G_T(\theta^{final}, \tau_x^{final}, s_h^{final}; T);$   

 $X' \leftarrow X + T';$   

Return:  $T'$ 

```

Figure 18: Algorithmic method for generating adversarial examples by Coe *et al.*

Algorithm 3 Attack Detection & Resilient State Estimation

Input: Observation $\{y_k\}_{k \geq 1} \in \mathbb{R}^q$; detection threshold η ; model parameters $A, B, C, \Sigma_v, \Sigma_w, \{u_k\}_{k \geq 0}$.
Output: Estimated values $\hat{x}_k, k \geq 1$

Initialize: $\hat{x}_0 \in \mathbb{R}^n$ and $P_0 \in \mathbb{R}^{n \times n}$.

- 1: **for** $k = 0, 1, 2, \dots$ **do**
- 2: Calculate $\hat{x}_{k+1|k}$ and $P_{k+1|k}$ using (4.5).
- 3: Apply χ^2 detector to
- 4: $y_{k+1} = (y_{k+1} - C\hat{x}_{k+1|k})^T (CP_{k+1|k}C^T + \Sigma_v)^{-1} (y_{k+1} - C\hat{x}_{k+1|k})$
- 5: **if** $y_{k+1} > \eta$ **then** $P_{k+1} = P_{k+1|k}$ and
- 6: $\hat{x}_{k+1} = \hat{x}_{k+1|k}$.
- 7: **else** Calculate P_{k+1} and \hat{x}_{k+1} using (4.6) and (4.7)
- 8: **end if**
- 9: Return: \hat{x}_{k+1} .
- 10: **end for**

Figure 19: Raj Gautam Dutta's (2018) attack detection method for adversarial attack generation

190 Results from LiDAR spoofing attacks show that the number of points increases robustness in machine
191 learning models [8]. As can be seen in the image below, the number of spoofed points decreases
192 robustness (as it results in an eighty-two percent success rate). When forty points are spoofed, fewer
193 spoofed points lead to an eighty-seven percent success rate and sixty spoofed points leads to a ninety
194 percent success rate [8]. LiDAR is used to make 3D representations of an area using data points
195 collected using laser beams. We can inject spoofed data points into a LiDAR by firing lasers into the
196 LiDAR's 2D camera's field of view [8]. As can be seen, the system can be fooled, and models should
197 be trained with spoofed data points to increase model robustness.

Targeted position	# Spoofed points		
	20	40	60
2-8 meters	87%	82%	90%

Figure 20: Number of spoofed points over the accuracy rate for LiDAR spoofing attacks [8]

198 3.2 Planning Module Attacks

199 The planning module can be attacked to fool a vehicle into making poor decisions. Poor decisions
200 making can prove fatal for an autonomous vehicle. We will investigate algorithms and data evaluating
201 such attacks.

202 **Planning Module Vulnerabilities** In the paper by Wan *et al.*, Simplified pseudo-code is analyzed
203 and tested to find vulnerabilities in the planning module.

```

1 # Iterate over the obstacle list
2 for (each obs : obstacle_list):
3     # Judge based on the lateral position of a veh.
4     # start_l, end_l are veh's left/right boundaries
5     if (obs.end_l < -2.5 or obs.start_l > 2.5):
6         continue
7     # The backward safe buffer
8     BackwardSafeBuffer ← 4.0f
9     # Check whether the veh is close
10    if (ego_start_s - obs.end_s < BackwardSafeBuffer):
11        IsClearChangeLane ← false

```

Figure 21: The authors iterate over the obstacle list and test to find obstacles with incorrect behavioural responses [10]

204 A vulnerability occurs when any attack target position is triggered [10]. As can be seen, triggering
205 objects can influence agent behaviour to cause agents to execute incorrect planning behaviour. The
206 violated PI is the violated planning invariant (which is the output behaviour given a set of vehicle,
207 obstacle, and pedestrian positions.

Vuln #	Driving Scenario	Software	Violated PI #	Attack-influenced Planning Behavior	Triggering Objects
V1	Lane following	Apollo 3.0/5.0	P11 (PI-C1)	Permanent stop	Static obstacle
V2	Lane changing	Apollo 3.0/5.0	P13 (PI-C2, 3)	Fail to change lane and never reach destination	Vehicle
V3	Lane borrow	Apollo 3.0/5.0	P14 (SP-PI-C1, 2)	Fail to borrow the lane and permanent stop	Vehicle or static obstacles in front of blocking vehicle
V4	Lane borrow	Apollo 3.0/5.0	P14 (PI-C1, 4, 5)	Fail to borrow the lane and permanent stop	Off-road static obstacle
V5	Intersection w/ traffic signal	Apollo 3.0/5.0	P16 (PI-C4)	Fail to pass intersection and permanent stop	Object in the way
V6	Intersection w/ traffic signal	Apollo 3.0/5.0	P16 (PI-C5)	Emergency stop possible to cause permanent stop and then fail to pass intersection	Pedestrian who is leaving the intersection
V7	Intersection w/ stop sign	Apollo 3.0/5.0	P15 (PI-C1)	Fail to pass intersection and permanent stop	Static bicycle off the road
V8	Lane following	Autoware	P11 (PI-C1)	Permanent stop	Static obstacle off the lane
V9	Lane following	Autoware	P11 (PI-C3)	Emergency stop	Moving vehicle off the lane

Figure 22: A list of driving scenarios violated planning invariants (PI's) based on a given triggering object which results in new attack-influenced planning behaviour [10]

208 **Adversarial Policy Attacks** The policy consists of a player playing against an opponent in a
 209 two-player Markov game. Adversaries are given unlimited black-box actions. Attackers must solve
 210 the given single-player Markov Decision Process [3]:

$$T_\alpha(s, a_\alpha) = T(s, a_\alpha, a_\nu) \quad \text{and} \quad R'_\alpha(s, a_\alpha, s') = R_\alpha(s, a_\alpha, a_\nu, s')$$

Figure 23: Single-Player Markov Decision Process formulation for the AI two-player zero-sum games presented [3].

211 The goal of an attacking agent is to find adversarial policies such that the sum of discounted rewards
 212 is maximized, as we can see below [3]:

$$\sum_{t=0}^{\infty} \gamma^t R_\alpha(s^{(t)}, a_\alpha^{(t)}, s^{(t+1)}), \quad \text{where } s^{(t+1)} \sim T_\alpha(s^{(t)}, a_\alpha^{(t)}) \text{ and } a_\alpha \sim \pi_\alpha(\cdot | s^{(t)})$$

Figure 24: Maximization for the sum of discounted rewards [3].

213 We note that the figure below shows the win rates against the medium victim in each AI-competitive
 214 environment [3]:

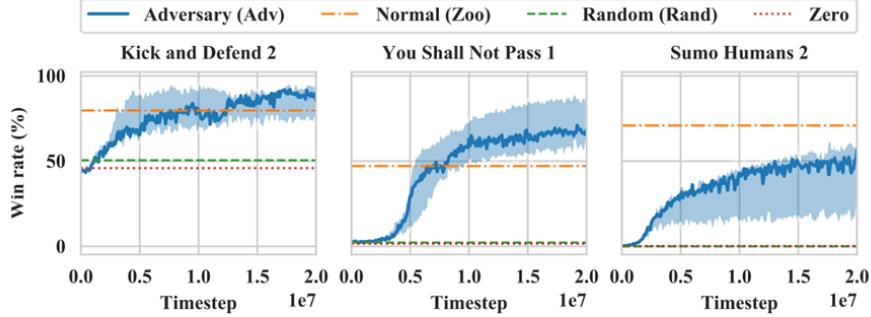


Figure 3: Win rates while training adversary Adv against the median victim in each environment (based on the difference between the win rate for Adv and Zoo). The adversary outperforms the Zoo baseline against the median victim in Kick and Defend and You Shall Not Pass, and is competitive on Sumo Humans. For full results, see figure 4 below or figure C.1 in the supplementary material.

Figure 25: Win-Rates while training adversary against the median victim in each zero-sum game [3]

215 Finally, the figure below shows the percentage of games won by an opponent given the policy used
 216 by the opponent [3]. Please note that zoo refers to the pre-trained baseline, zero refers to the lifeless
 217 policy with zero control and refers to the random action policy, and adv refers to the adversarial
 218 policy:

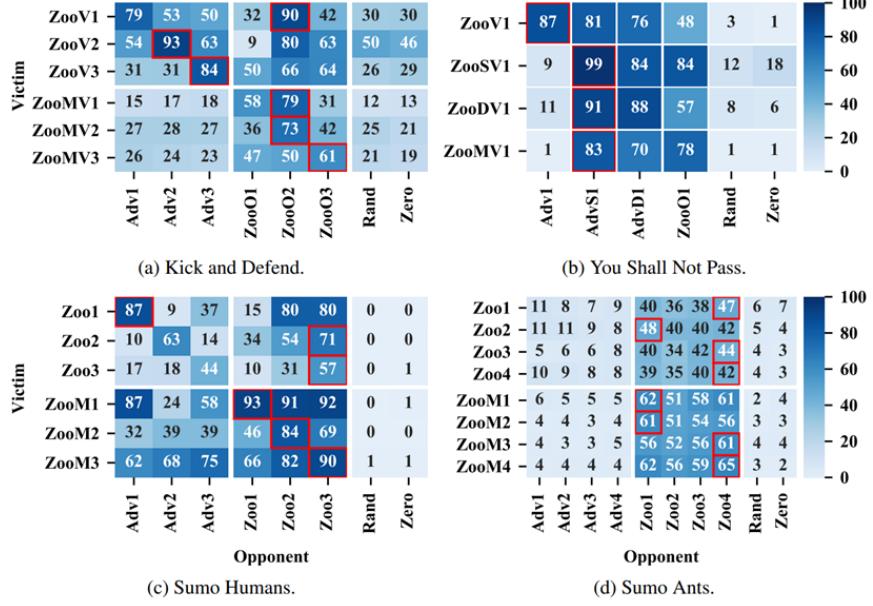


Figure 26: As can be seen, carefully calculated mistakes can be used to fool AI models to eliminate good agents from competitive tasks [3]

219 **Planning Around Physical Obstacles** The table below shows the selected variables for attack
 220 injection. Correlated variables indicate a connection between properties measured by physical sensors
 221 and expected vehicle behaviour. If a vehicle’s machine learning based planning-module is trained
 222 using sensor data, then correlations between sensor values are mapped within the machine learning
 223 neural network. Changing a sensor value via an adversarial attack can trigger invalid inferences from
 224 the neural network since it expects sensor values to be correlated for given planning-module-based
 225 decisions. Note that FPR denotes the false positive rate for spoofing detection, and correction denotes
 226 the sensor correlation (a correlation is significant if it has a value greater than a half) [11].

Variable 1	Variable 2	Correlation	FPR
position.y	linear_velocity.y	-0.741	8.77%
position.z	orientation.qz	0.801	12.65%
orientation.qy	euler_angles.y	0.877	3.18%
position_std.dev.x	orientation_std.dev.x	0.845	2.22%
orientation.qw	euler_angles.y	-0.801	6.54%
linear_velocity.x	euler_angles.z	-0.964	2.55%
linear_acceleration.x	linear_acceleration_vrf.x	0.899	4.66%
angular_velocity.z	angular_velocity_vrf.z	0.999	0.3%
orientation_std.dev.z	linear_velocity_std.dev.z	0.699	9.20%
linear_velocity_std.dev.x	linear_velocity_std.dev.y	0.915	0.83%

Figure 27: Correlation values between vehicular sensor variables [11]

227 Figure (a) below shows 184 variable pairwise correlations in a matrix. Figure (b) shows we can cause
 228 errors in run-time correlation, which can compromise a vehicle’s planning module. Therefore, there
 229 is empirical evidence of neural network models being compromised when correlated variable pairs
 230 have a variable which is attacked [11].

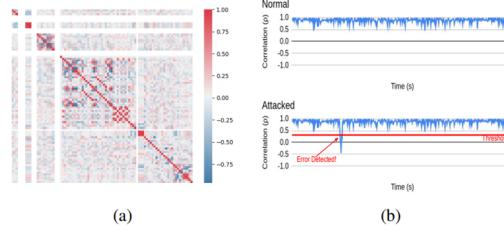


Figure 28: Correlation between variables (a) and resulting compromised vehicular planning (b) during adversarial sensor attacks [11]

231 3.3 Testing and Evaluation of Attacks

232 We explore various methods centred around the testing and evaluation of attack methods for au-
233 tonomous vehicles. Understanding what makes a vehicle vulnerable to attack, how to plan around
234 attacks, test vulnerabilities, and debug vulnerabilities is essential to understand attacks on autonomous
235 vehicle models.

236 **Plan, Test, and Debugging** Robustness values close to zero represent cases where the vehicle
237 collides with an object at low speeds. The probability distribution is a truncated normal distribution.
238 Global Uniform Random Search (Global UR) randomly selects test cases for simulation and checks
239 if traces satisfy constraints [4]. Covering Arrays and Uniform Random Search (CA+UR) creates a
240 list of test cases and evaluations against specifications [4]. Covering Arrays and Simulated Annealing
241 (CA+SA) combines covering arrays to evaluate discrete variables and simulates annealing, so a
242 function is used to guide the search of continuous variables [4].

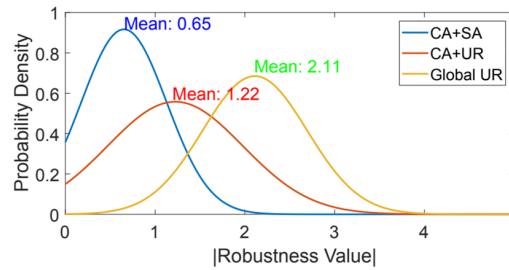


Figure 29: We can see that the CA+SA system shows the best results given its probability distribution clustering close to zero on the x-axes

243 **Defences for Neural Networks** Algorithms can generate adversarial examples to reduce confidence
244 in classification, misclassify images, make targeted misclassifications, and cause source-target
245 misclassifications [7].

```
Algorithm 1: Generating adversarial samples
Input: X, Y*, F, T, θ
1: X' ← X
2: ε ← ε(X)
3: while F(X') ≠ Y* and ||δX|| < T do
4:   Compute forward derivative  $\nabla F(X, Y)$ 
5:    $\delta X = \text{sign}(\nabla F(X, Y)) \cdot \epsilon$ 
6:   Modify  $X' \leftarrow X' + \theta \cdot \text{sign}(\nabla F(X, Y)) \cdot \delta X$ 
7:    $\theta \leftarrow \theta + \eta$ 
8: end while
9: return X'
```

```
Algorithm 2: Generating adversarial samples for FGSM
Input: X, Y*, F, T, θ
1: X' ← X
2: ε ← ε(X)
3: while F(X') ≠ Y* and ||δX|| < T do
4:   Compute forward derivative  $\nabla F(X, Y)$ 
5:    $\delta X = \text{sign}(\nabla F(X, Y)) \cdot \epsilon$ 
6:   Modify  $X' \leftarrow X' + \theta \cdot \text{sign}(\nabla F(X, Y)) \cdot \delta X$ 
7:    $\theta \leftarrow \theta + \eta$ 
8: end while
9: return X'
```

```
Algorithm 3: Increasing pixel intensities saliency map
Input:  $\nabla F(X), \Gamma, \ell$ 
1: for each pair  $(p, q) \in \Gamma$  do
2:    $\alpha = \sum_{i \in \text{source class}} \frac{\partial \nabla F_i}{\partial X_p}$ 
3:    $\beta = \sum_{i \in \text{target class}} \sum_{j \in \Gamma \setminus p} \frac{\partial \nabla F_i}{\partial X_j}$ 
4:   if  $\alpha > 0$  and  $\beta < 0$  and  $-\alpha < \beta > \max$  then
5:      $p_1, p_2 \leftarrow \text{argmax}_{i \in \Gamma} \nabla F_i(p, q, \ell, Y^*)$ 
6:      $\delta X_p \leftarrow \text{sign}(\nabla F(p_1, p_2, \ell, Y^*)) \cdot \epsilon$ 
7:     Remove  $p$  from  $\Gamma$  if  $p_1$  and  $p_2$  are not in  $\Gamma$ 
8:      $\ell \leftarrow \ell + \eta$ 
9:   end if
10: end for
11: return  $p_1, p_2$ 
```

Figure 30: Algorithms to generate and test adversarial examples [7]

246 Adversarial examples can be used to create quantitative metrics for sampling robustness.

Source set of 10,000 original samples	Adversarial samples successfully misclassified	Average distortion	
		All adversarial samples	Successful adversarial samples
Training	97.05%	4.45%	4.03%
Validation	97.19%	4.41%	4.01%
Test	97.05%	4.45%	4.03%

Figure 31: Empirical results of adversarial testing [7]

247 4 Remarks

248 As previous authors have significantly elaborated on the data behind attacks, the remainder of the
 249 paper will focus on a theoretical conclusion of the meaning behind the attacks.

250 4.1 Perceptron Attacks

251 Perceptron models appear vulnerable to attacks since they lack a proper grounding in the meaning
 252 behind the object shapes they detect. Baloor *et al.* showed in their work that physical attacks such
 253 as placing lines on objects could be used to fool neural network models into making environmental
 254 misclassifications. Algorithms to detect attacks and train machine learning models under adversarial
 255 attacks have been proposed to prevent them, but they do not resolve the problems associated with
 256 perceptron models not understanding underlying properties to distinguish objects. Determining
 257 if a perceptron learning algorithm converges on a correct shape without defining the shape leaves
 258 it vulnerable to attacks based on point and line-based distortion. Vehicles driving on roads are
 259 vulnerable to attacks since they learn by approximating shape outlines, meaning they cannot become
 260 perfectly robust when using neural network models. Lines and patches added to roads and stop signs
 261 cause object misclassification by creating incorrect inferences for polygons within the perceptron
 262 model.

263 4.2 Planning Attacks

264 Planning attacks seem generalizable because they are vulnerable at points between vehicle goal
 265 setting, vehicle action planning, and vehicle action completion. Vehicles can incorrectly plan goals
 266 based on invalid perceptron input leading to inaccurate actions being made on the road. Furthermore,
 267 adversarial drivers can interfere with planning systems since they can provide otherwise unexpected
 268 input to driving vehicles, leading to undefined driving behaviours. Planning networks are trained by
 269 determining optimal moves given regular input, meaning training networks with a global set of rules
 270 to safeguard against adversarial drivers could be beneficial. Multiple network models for different
 271 aspects of the road, like distance from the curb, can be used to create a priority of perception-based
 272 plans. To ignore obstacles along the road, multiple sets of sensors to determine objects outside the
 273 road to be ignored can be added.

274 5 Conclusion

275 There are many ways to fool a planning module for a machine learning system. Determining how
 276 the planning module and perceptron module are interrelated and determining gives excellent insight
 277 into the meaning behind machine learning models. Determining how we test planning modules given
 278 the large number of choices they are presented with is an essential question for future researchers.
 279 Finally, evaluating if our methods are trustworthy and how we can collect more data on adversarial
 280 vehicle attacks is an important step forward for future research.

281 **References**

- 282 [1] A. Boloor, K. Garimella, X. He, C. Gill, Y. Vorobeychik, and X. Zhang, “Attacking vision-based perception
283 in end-to-end autonomous driving models,” Journal of Systems Architecture, vol. 110, p. 101766, 2020.
- 284 [2] A. Chowdhury, G. Karmakar, J. Kamruzzaman, A. Jolfaei, and R. Das, “Attacks on self-driving cars and their
285 countermeasures: A survey,” IEEE Access, vol. 8, pp. 207308–207342, 2020.
- 286 [3] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine, and S. Russell, “Adversarial Policies: Attacking Deep
287 Reinforcement Learning,” 2019.
- 288 [4] C. E. Tuncali, G. Fainekos, H. Ito, and J. Kapinski, “Simulation-based adversarial test generation for
289 autonomous vehicles with machine learning components,” 2018 IEEE Intelligent Vehicles Symposium (IV),
290 2018.
- 291 [5] J. Zhang, Y. Lou, J. Wang, K. Wu, K. Lu, and X. Jia, “Evaluating adversarial attacks on driving safety in
292 vision-based Autonomous Vehicles,” IEEE Internet of Things Journal, vol. 9, no. 5, pp. 3443–3456, 2022.
- 293 [6] R. G. Dutta and Y. Jin, “Security of autonomous systems under physical attacks: With application to
294 self-driving cars,” dissertation, STARS, Orlando, FL, 2018.
- 295 [7] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, “The Limitations of Deep
296 Learning in Adversarial Settings,” arXiv, 2015.
- 297 [8] Y. Cao, C. Xiao, B. Cyr, Y. Zhou, W. Park, S. Rampazzi, Q. A. Chen, K. Fu, and Z. M. Mao, “Adversarial
298 sensor attack on Lidar-based perception in autonomous driving,” Proceedings of the 2019 ACM SIGSAC
299 Conference on Computer and Communications Security, 2019.
- 300 [9] Z. El-Rewini, K. Sadatsharan, N. Sugunaraj, D. F. Selvaraj, S. J. Plathottam, and P. Ranganathan, “Cyberse-
301 curity attacks in vehicular sensors,” IEEE Sensors Journal, vol. 20, no. 22, pp. 13752–13767, 2020.
- 302 [10] Z. Wan, J. Shen, J. Chuang, X. Xia, J. Garcia, J. Ma, and Q. A. Chen, “Too afraid to drive: System-
303 atic discovery of semantic DOS vulnerability in autonomous driving planning under physical-world attacks,”
304 Proceedings 2022 Network and Distributed System Security Symposium, 2022.
- 305 [11] Z. Zhang, S. S. V. Singapuram, Q. Zhang, D. K. Hong, B. Nguyen, Z. M. Mao, S. Mahlke, and Q. A. Chen,
306 “AVMaestro: A Centralized Policy Enforcement Framework for Safe Autonomous-driving Environments.”

307 **Checklist**

308 See the questions and answers below:

- 309 1. For all authors...
 - 310 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
311 contributions and scope? **[Yes]**
 - 312 (b) Did you describe the limitations of your work? **[Yes]**
 - 313 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - 314 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
315 them? **[N/A]**
- 316 2. If you are including theoretical results...
 - 317 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 318 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 319 3. If you ran experiments...
 - 320 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
321 mental results (either in the supplemental material or as a URL)? **[N/A]**
 - 322 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
323 were chosen)? **[N/A]**
 - 324 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
325 ments multiple times)? **[N/A]**
 - 326 (d) Did you include the total amount of compute and the type of resources used (e.g., type
327 of GPUs, internal cluster, or cloud provider)? **[N/A]**
- 328 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 329 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - 330 (b) Did you mention the license of the assets? **[Yes]**
 - 331 (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
 - 332 (d) Did you discuss whether and how consent was obtained from people whose data you're
333 using/curating? **[N/A]**
 - 334 (e) Did you discuss whether the data you are using/curating contains personally identifiable
335 information or offensive content? **[N/A]**
- 337 5. If you used crowdsourcing or conducted research with human subjects...
 - 338 (a) Did you include the full text of instructions given to participants and screenshots, if
339 applicable? **[TODO]**
 - 340 (b) Did you describe any potential participant risks, with links to Institutional Review
341 Board (IRB) approvals, if applicable? **[TODO]**
 - 342 (c) Did you include the estimated hourly wage paid to participants and the total amount
343 spent on participant compensation? **[TODO]**