

# Effects of Learning-Based Action-Space Attacks on Autonomous Driving Agents

Yuting Wu, Xin Lou<sup>† ‡</sup>, Pengfei Zhou<sup>††</sup>, Rui Tan, Zbigniew T. Kalbarczyk\*, Ravishankar K. Iyer\*

Nanyang Technological University, Singapore

<sup>†</sup> Singapore Institute of Technology, <sup>‡</sup> Illinois at Singapore

<sup>††</sup> University of Pittsburgh, USA

\*University of Illinois at Urbana-Champaign, USA

## ABSTRACT

Vehicle cybernation with increasing use of information and communication technologies faces cybersecurity threats. This extended abstract studies action-space attacks on autonomous driving agents that make decisions using either a traditional modular processing pipeline or the recently proposed end-to-end model obtained via deep reinforcement learning (DRL). The action-space attacks alter the actuation signal and pose direct risks to the vehicle's behavior. We formulate the attack construction as a DRL problem based on the input from either an extra camera or inertial measurement unit deployed. Attacks are designed to lurk until a safety-critical moment arises (e.g. lane changing or overtaking), with the goal of causing a side collision upon activation. Our results demonstrate that the modular processing pipeline is more resilient than the DRL-based agent, due to the former's main focus of trajectory following. We further investigate two enhancement methods: adversarial training through fine-tuning and progressive neural networks, gaining an essential understanding of their pros and cons.

## 1 INTRODUCTION

Recent rapid growth in autonomous driving (AD) has brought research attention to its cybersecurity concerns. Rising autonomy results in more sensors and connectivity, thereby expanding potential attack targets in AD. Among miscellaneous possible attack mount points, targeting the actuation of a vehicle is appealing to the attacker. Adversaries can bypass potential defense mechanisms and directly affect the vehicle's state. However, action-space attacks, also referred to as actuator attacks, have gained limited attention in the context of AD. Most recent studies on action-space attacks in AD rely on model-based approaches that either require in-vehicle data for the current system state [6] or vehicle's kinematics and structure [2], resulting in a demanding form of white-box attacks. Meanwhile, attacks in the black-box setting have been mostly studied in simulation environments like OpenAI Gym [3] and Mathworks [5], which are not representative of real-world

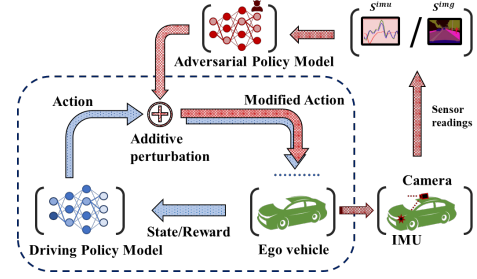


Figure 1: Overview of the DRL-based action-space attack.

driving conditions. Additionally, these studies usually concentrate on a single type of driving system without contrasting the impact of the attack across various AD designs. In light of this, we study the action-space attack on two major types of driving agents: 1) the traditional modular processing pipeline and 2) the end-to-end policy model trained via DRL. To ensure the realism of the attack, we assume the attacker has no access to (i) the driving agent's internal, and (ii) the driving agent's sensor readings. Both driving agents are formulated based on a trajectory-following task while adopting different design methodologies. We hypothesize that the design differences between the two agents will lead to distinct characteristics in responding to action-space attacks. We further apply adversarial training to enhance the DRL-based end-to-end driving agent with two variants, fine-tuning and progressive neural networks (PNN).

## 2 METHODOLOGY

We treat the entire driving system as a black box, utilizing DRL to investigate safety-critical moments and to learn how to introduce disturbances, as depicted in Fig. 1. In our study, the attacker utilizes either an extra camera or an inertial measurement unit (IMU) to identify safety-critical moments in the driving system. The former provides adequate information while its installation demands a wide field of view, which may attract attention from humans. The latter provides a less informative inertia trace but can be concealed within the vehicle, making them nearly undetectable. To explore IMU's potential utilization by attackers, we proposed a 'learning-from-teacher' structure to transfer the attack policy from obtained camera-based attacks to IMU-based attacks. The action-space attack injects additive perturbations into the steering angle of the ego vehicle at safety-critical moments (i.e., lane changing and overtaking), aiming to create a side collision with another vehicle on road. The attack is subjected to an attack budget that characterizes the actuation system's logistic constraint (i.e., the maximum allowed adjustment value per actuation step) or the desired degree of attack stealthiness. We conduct two variants of adversarial training to

<sup>1</sup>This project is supported by the National Research Foundation, Singapore, and the National University of Singapore through its National Satellite of Excellence in Trustworthy Software Systems (NSOE-TSS) office under the Trustworthy Computing for Secure Smart Nation Grant (TCSSNG) award no. NSOE-TSS2020-01, and in part by the National Research Foundation, Prime Minister's Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) program.

<sup>2</sup>Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s). ICCPS '23, May 9–12, 2023, San Antonio, TX, USA © 2023 Copyright held by the owner/author(s). ACM ISBN 979-8-4007-0036-1/23/05. <https://doi.org/10.1145/3576841.3589615>

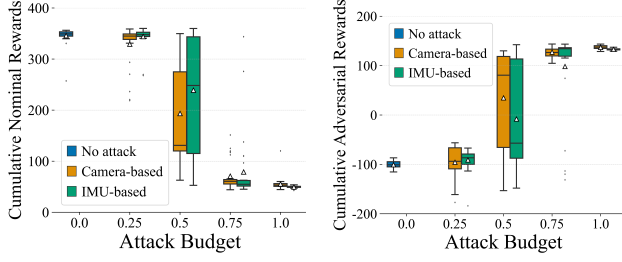


Figure 2: Evaluation under various attack budgets.

enhance the end-to-end driving policy: fine-tuning and PNN. For the former one, we control the ratio of selecting zero attack budget (i.e., no attack) to prevent overfitting to adversarial cases and forgetting the nominal driving pattern. For the latter one, we adjust the threshold for switching between either the original column or the adversarially trained column based on the attack budget.

### 3 EXPERIMENT RESULTS AND DISCUSSION

We conduct experiments in CARLA 0.9.11 [1], with CARLA Autopilot as the modular driving pipeline. We use DRL with a path planner in the reward design to construct the end-to-end driving agent. Both camera-based and IMU-based attacks successfully learned a policy that perturbs the steering angle of the victim vehicle at lane-changing and overtaking moments, resulting in a side collision with other vehicles on road. Comparative results of the attack effectiveness under various attack budgets are given in Fig. 2. Nominal driving rewards assess driving performance based on speed and traveling distance, and adherence to the planned path. Adversarial rewards assess attack accuracy and success rate. These results suggest a trade-off between the accuracy and stealthiness of the attack. When the attacker has direct camera observation, it is easier for the attacker to time the hijacking of the ego vehicle and remain silent when the attack conditions are not met. However, when the attacker only has indirect observations via stealthily installed IMU, identification of safety-critical moments to produce a desired attack impact becomes more challenging. We assess the robustness of two driving agents based on the correlation between steering deviation from the predetermined path and the attack effort, which is averaged injected perturbations applied over each attack attempt. As shown in Fig. 3, compared with the end-to-end driving agent, the modular driving agent can maintain more minor tracking errors in the trajectory following task when the attack effort is low. The superior performance of the modular driving agent can be attributed to the inclusion of a proportional–integral–derivative (PID) controller in its design. The PID controller calculates the throttle, brake, and steering output needed to keep the vehicle on the planned path. Once an error in speed heading is observed, the PID controller adjusts the actuator value instantly. In contrast, the end-to-end driving agent is trained to optimize a linear combination of multiple goals, resulting in its tendency of trading precision for faster speeds. We further explored the performance of the enhanced agents. Each enhancement involves two variants. As shown in Fig. 4, the enhanced agent with fine-tuning serves better in the presence of attacks in terms of trajectory following accuracy. Yet, the catastrophic forgetting problems cause degraded driving performance

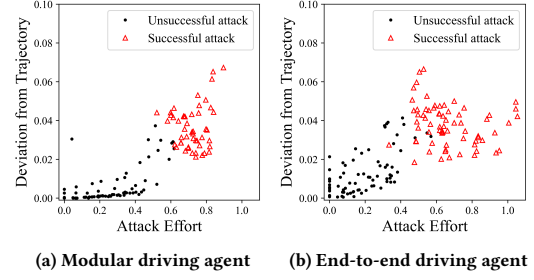


Figure 3: Evaluation of different driving agents under attacks.

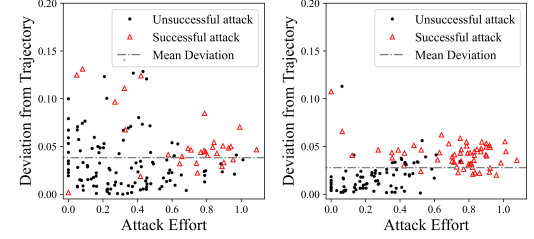


Figure 4: Enhanced agent with fine-tuning.

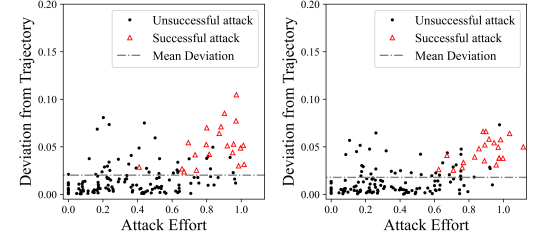


Figure 5: Enhanced agent with progressive neural networks.

with larger deviations from the planned path, even without significant attacks. To avoid the forgetting problem, we further extend the driving agent into PNN. As shown in Fig. 5, when faced with attacks, it outperforms the former while maintaining its nominal driving behavior by switching between a nominal and adversarially trained driving policy. However, it has limitations since it requires the identification of attacks to switch between policies. Although action-space attacks are rare, they cannot be ignored. Therefore, our results suggest that a simplex driving agent [4] capable of switching between the enhanced driving policy model and the nominal driving agent when attacks are detected is desirable.

### REFERENCES

- [1] A. Dosovitskiy, G. Ros, F. Codevilla, A. M. López, and V. Koltun. 2017. CARLA: An Open Urban Driving Simulator. In *CoRL*, Vol. 78. PMLR, 16.
- [2] S. Jha, S. Banerjee, T. Tsai, S. KS Hari, M. B Sullivan, Z. T Kalbarczyk, S. W Keckler, and R. K Iyer. 2019. ML-based fault injection for autonomous vehicles: A case for bayesian fault injection. In *2019 49th annual IEEE/IFIP DSN*. IEEE, 112–124.
- [3] X. Y. Lee, Y. Esfandiari, K. L. Tan, and S. Sarkar. 2021. Query-based targeted action-space adversarial policies on deep reinforcement learning agents. In *ICCPs*. ACM, 87–97.
- [4] S. Mohan, S. Bak, E. Betti, H. Yun, L. Sha, and M. Caccamo. 2013. S3A: secure system simplex architecture for enhanced security and robustness of cyber-physical systems. In *HiCoNS*. ACM, 65–74.
- [5] M. Moradi, B. James Oakes, M. Saraoglu, A. Morozov, K. Janschek, and J. Denil. 2020. Exploring fault parameter space using reinforcement learning-based fault injection. In *2020 50th Annual IEEE/IFIP DSN-W*. IEEE, 102–109.
- [6] X. Zhou, A. Schmedding, et al. 2022. Strategic safety-critical attacks against an advanced driver assistance system. In *2022 52nd Annual IEEE/IFIP ICDNS*. IEEE.