

# Revisiting Adversarial Perception Attacks and Defense Methods on Autonomous Driving Systems

Cheng Chen, Yuhong Wang, Nafis S Munir, Xiangwei Zhou, Xugui Zhou

Louisiana State University, Baton Rouge, LA 70803 {cchen72, ywan248, nmunir1, xwzhou, xuguizhou}@lsu.edu

**Abstract**—Autonomous driving systems (ADS) increasingly rely on deep learning-based perception models, which remain vulnerable to adversarial attacks. In this paper, we revisit adversarial attacks and defense methods, focusing on road sign recognition and lead object detection and prediction (e.g., relative distance). Using a Level-2 production ADS, OpenPilot by Comma.ai, and the widely adopted YOLO model, we systematically examine the impact of adversarial perturbations and assess defense techniques, including adversarial training, image processing, contrastive learning, and diffusion models. Our experiments highlight both the strengths and limitations of these methods in mitigating complex attacks. Through targeted evaluations of model robustness, we aim to provide deeper insights into the vulnerabilities of ADS perception systems and contribute guidance for developing more resilient defense strategies. [Code Available at [https://github.com/DepCPS/revisiting\\_adversarial\\_ADS](https://github.com/DepCPS/revisiting_adversarial_ADS)]

**Index Terms**—ADS, adversarial attack, defense, autonomous vehicle.

## I. INTRODUCTION

Autonomous driving systems (ADS) are at the forefront of revolutionizing modern mobility, transforming our roads into safer and more efficient pathways for the future [1]. Every major car manufacturer has embraced this wave by developing their own variants of autonomous driving or assisted driving solutions (e.g., Tesla’s Full Self-Driving (FSD) [2]), underscoring both the competitive and innovative spirit in the industry. The key to ADS lies in accurate and powerful environment perception modules, primarily enabled through deep learning-based perception models [3]. However, these models remain highly susceptible to adversarial attacks, deliberate perturbations designed to deceive perception systems and trigger dangerous decisions or system failures [4]–[6].

The challenges are substantial. As adversarial techniques grow more sophisticated, existing defense strategies struggle to keep pace, often remaining confined to specially trained models and curated datasets. This disconnect between controlled testing environments and the unpredictable conditions of real-world driving underscores the urgent need for comprehensive and realistic investigations into ADS vulnerabilities and the development of robust, scalable defense mechanisms.

Reliable traffic sign recognition is vital for maintaining situational awareness and ensuring compliance with road rules. Disruptions in this process can result in improper speed regulation or navigation errors. Moreover, accurate distance estimation is essential for tasks such as collision avoidance and autonomous navigation. Adversarial perturbations that distort distance measurements can undermine the safety margins required for secure driving.

Motivated by these challenges, this paper systematically revisits adversarial attacks and defense strategies, integrating them into the open-source commercial ADS framework. We focus on two critical perception tasks: traffic sign classification and distance regression, using YOLOv8 [7] and OpenPilot’s Supercombo model as case studies to evaluate these threats comprehensively [8]. We experimentally examine several adversarial techniques, including Gaussian Noise [9], Fast Gradient Sign Method (FGSM) [10], Auto Projected Gradient Descent (Auto-PGD) [11], Simple Black-box Attack (SimBA) [12], Robust Physical Perturbations (RP<sub>2</sub>) [13], and CAP-attacks [14]. To mitigate these vulnerabilities, we evaluate a range of defense methods, including adversarial training, image processing techniques, contrastive learning, and diffusion-based approaches. Through this comprehensive analysis, our research aims to advance the understanding of adversarial risks in ADS and support the development of more resilient and trustworthy autonomous driving systems.

Our research aims to deepen the understanding of adversarial vulnerabilities in ADS perception models and offer actionable insights into effective defense strategies, ultimately advancing the safety and resilience of autonomous driving systems. Our main contributions include:

- We systematically revisit and compare the threats and limitations of different adversarial attack algorithms.
- We evaluate existing defense mechanisms in conjunction with different attack types to assess their effectiveness in ADS contexts.
- We propose directions and perspectives for advancing research on adversarial attacks and defenses in autonomous driving.

## II. BACKGROUND

### A. Autonomous Driving System

Autonomous vehicles rely on sensors, actuators, advanced algorithms, often driven by machine learning, and high-performance processors to execute software. They construct and maintain a perceptual map of their surroundings using diverse sensors embedded throughout the vehicle [15]. For instance, radar sensors monitor the position of nearby vehicles [16], while cameras detect traffic lights, read road signs, track other vehicles, and identify drivable paths [17]–[19]. LiDAR sensors emit light pulses to measure distances, detect road boundaries, and recognize lane markings [20]. These multimodal inputs are processed by complex models that plan

driving paths and issue commands to actuators responsible for acceleration, braking, and steering.

### B. Adversarial Attack

Adversarial attacks are deliberate manipulations designed to mislead machine learning models by introducing subtle, carefully crafted perturbations to input data [21]. Although often imperceptible to humans, these perturbations can significantly degrade model performance, leading to incorrect or unsafe outputs. In safety-critical domains like autonomous driving, adversarial attacks present a substantial risk, potentially triggering hazardous behaviors or system failures [22]. Identifying and mitigating these vulnerabilities is essential for ensuring the robustness and dependability of autonomous systems.

### C. Adversarial Defense

Adversarial defense strategies aim to enhance the robustness of machine learning models against adversarial perturbations [23]. Common approaches include adversarial training, where models are trained using adversarial examples; input preprocessing techniques such as image denoising and geometric transformations; and advanced representation learning methods like contrastive learning and diffusion-based generative models. Developing effective defenses is especially critical in autonomous driving to preserve the reliability and safety of perception systems under adversarial conditions.

## III. PERCEPTION ATTACKS

ADS heavily relies on perception modules, which are typically built using deep neural networks to process data from various sensors. However, these perception modules are highly susceptible to adversarial attacks. In this section, we introduce several common adversarial attacks that have been extensively studied and pose significant threats to ADS.

### A. Gaussian Noise

Original images, during acquisition and transmission, are often affected by noise, which degrades image quality, blurs features, and complicates analysis. Gaussian noise, characterized by a probability density that follows a Gaussian distribution, represents one of the simplest adversarial methods. This attack introduces random perturbations into the input data by adding noise sampled from a Gaussian distribution:

$$x_{adv} = x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2). \quad (1)$$

where  $x_{adv}$  is the adversarial example,  $x$  is the original input, and  $\epsilon$  is additive noise drawn from a zero-mean normal distribution with variance  $\sigma^2$ . Although Gaussian noise is not specifically optimized against the model, it can notably degrade performance, particularly in environments with sensor uncertainties such as nighttime driving, fog, or rain, thereby exposing perception model vulnerabilities.

### B. Fast Gradient Sign Method

The FGSM, proposed by Goodfellow et al. [10], is a white-box attack that generates adversarial examples using gradient information from the neural network. By adding small, directionally aligned perturbations to the input, the model is tricked into incorrect predictions:

$$x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, y)). \quad (2)$$

where  $J(\theta, x, y)$  is the model's loss function with respect to the parameters  $\theta$ , input  $x$ , and true label  $y$ . The gradient sign indicates the direction in input space that maximally increases the loss. Due to its computational efficiency, FGSM is widely used for preliminary robustness assessments.

### C. Auto Projected Gradient Descent

Auto-PGD [11] extends gradient-based attacks by applying iterative updates and adaptive step sizes to improve perturbations. Each iteration calculates a perturbation step and projects it back into a feasible region, maintaining imperceptibility:

$$x_{adv}^{t+1} = \text{Proj}_{x+S} \left( x_{adv}^t + \alpha \cdot \text{sign} \left( \nabla_{x_{adv}^t} J(\theta, x_{adv}^t, y) \right) \right). \quad (3)$$

In this formulation,  $S$  defines the permissible perturbation space,  $\alpha$  is the step size, and Proj ensures the perturbation remains within a perceptual bound. Compared with simpler methods like FGSM, Auto-PGD produces much stronger adversarial examples, making it a powerful tool for evaluating the resilience of high-performance ADS perception systems.

### D. Simple Black-box Attack

SimBA, introduced by [12], is a query-efficient black-box adversarial attack that operates without gradient estimation. It iteratively refines an adversarial perturbation  $\delta$  through random sampling from an orthonormal basis (e.g., pixel or DCT frequency basis). In each step, a basis vector  $\mathbf{q}$  is chosen, and  $\delta$  is updated by adding or subtracting a small step  $\epsilon$  along  $\mathbf{q}$ . The update direction depends on output probabilities: in untargeted attacks, the direction that reduces the true class probability  $p(y_{\text{true}} | \mathbf{x} + \delta)$  is selected; in targeted attacks, the direction that increases the target class probability  $p(y_{\text{target}} | \mathbf{x} + \delta)$  is chosen.

The cumulative perturbation after  $T$  steps is bounded by:

$$\|\delta_T\|_2^2 \leq T\epsilon^2. \quad (4)$$

This bound ensures controlled perturbation growth. Despite its simplicity, SimBA's orthonormal sampling strategy enables high query efficiency, making it a strong baseline among black-box attacks.

### E. Optimization-based Methods

1) *Robust Physical Perturbations:* To assess deep neural network (DNN) vulnerabilities under realistic conditions, RP<sub>2</sub> [13] generates physical-world adversarial perturbations that remain effective across diverse environments, varying viewpoints, lighting, distances, and sensor noise. Given a classifier  $f\theta(x)$  and true label  $y$ , the goal is to find a perturbation  $\delta$  such

that the perturbed input  $x' = x + \delta$  is classified as a specific target label  $y^* \neq y$ , i.e.,

$$f_\theta(x + \delta) = y^*. \quad (5)$$

The optimization problem solved by  $\text{RP}_2$  is formulated as:

$$\arg \min_{\delta} \lambda \|\mathbf{M}_x \cdot \delta\|_p + \text{NPS} \\ + \mathbb{E}_{x_i \sim \mathcal{X}_V} [J(f_\theta(x_i + T_i(\mathbf{M}_x \cdot \delta)), y^*)]. \quad (6)$$

Here,  $\|\mathbf{M}_x \cdot \delta\|_p$  constrains the perturbation within the object's surface (e.g., a stop sign) using a binary mask  $\mathbf{M}_x$ , and NPS (Non-Printability Score) penalizes non-reproducible colors.  $\mathcal{X}_V$  represents images under varying conditions, and  $T_i$  denotes transformations (e.g., rotation, scaling). The function  $J(\cdot)$  is typically cross-entropy loss. These perturbations highlight the real-world feasibility of adversarial threats and their implications for autonomous driving safety.

2) *CAP-Attack*: CAP-Attack [14] introduces a runtime adversarial patch generation technique for DNN-based Adaptive Cruise Control (ACC) systems. Unlike offline adversarial generation, this approach performs real-time pixel-level tuning on the front vehicle's image region to mislead distance prediction models:

$$\min_{\Delta_t} \sum_{d \in RD_t} -\nabla g(d, \theta) + \lambda \|\Delta_t\|_p. \quad (7)$$

An attribution mechanism pinpoints regions most sensitive to predictions and disturbances  $\Delta_t$  are confined to the bounding box of the front vehicle to reduce computation and enhance stealth. The algorithm inherits and adapts perturbations frame-by-frame, adjusting based on the displacement and size of the vehicle, ensuring both temporal coherence and concealment. This continuous optimization enhances attack effectiveness while making detection by traditional safety systems more difficult. Here,  $\lambda$  controls the trade-off between the impact and invisibility of the patch, and  $\|\Delta_t\|_p$  regularizes the perturbation magnitude for stealthiness.

#### IV. DEFENSE METHODS

The robustness of DNNs in ADS is critical, yet adversarial attacks threaten their reliability and safety. Developing effective defense strategies to bolster adversarial robustness is therefore essential; this section outlines several widely used adversarial defense methods.

##### A. Image Processing

Applying simple and consistent preprocessing techniques to images before feeding them into the model is a straightforward and effective strategy. In this subsection, we introduce three such image processing techniques:

Median Blurring [24] applies median filtering to suppress adversarial noise by replacing each pixel with the median of its neighborhood, effectively preserving edges while mitigating perturbations.

Bit-depth Reduction [24] reduces the precision of pixel values, thereby decreasing the effectiveness of subtle perturbations and enhancing robustness against adversarial inputs.

Randomization [25] introduces randomness, such as random resizing, padding, or noise injection, during preprocessing, which disrupts adversarial perturbations and hinders consistent model exploitation.

##### B. Adversarial Training

Adversarial training [26] is a widely used and powerful defense technique that improves model robustness by solving a min-max optimization problem:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\|\delta\| \leq \epsilon} \mathcal{L}(f_\theta(x + \delta), y) \right]. \quad (8)$$

The inner maximization, which identifies the worst-case perturbation  $\delta$  within a bounded region defined by  $\epsilon$ ) finds the most adversarial example, while the outer minimization updates the model parameters  $\theta$  to reduce the loss  $\mathcal{L}$  on these perturbed inputs. By incorporating adversarial examples into the training process, the model learns locally stable decision boundaries. Despite its effectiveness, especially against white-box and some black-box attacks, adversarial training generally involves greater computational overhead and may also lead to reduced performance on clean inputs.

##### C. Diffusion Model

Diffusion models have demonstrated impressive capabilities in generating data that resembles real-world distributions. We leverage this property to repair adversarially attacked images. We employ DiffPIR [27], a restoration framework that uses a pre-trained diffusion model as a powerful generative prior to remove adversarial or noisy perturbations. Instead of training a separate denoising model, DiffPIR harnesses the score function of the diffusion process to guide image restoration while preserving fine visual details.

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \cdot \arg \min_{\mathbf{x}} \|\mathbf{y} - H(\mathbf{x})\|^2 + \rho_t \|\mathbf{x} - \mathbf{x}_0^{(t)}\|^2 \\ + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \left( \sqrt{1 - \zeta} \cdot \hat{\epsilon} + \sqrt{\zeta} \cdot \epsilon_t \right). \quad (9)$$

Restoration is performed via two lightweight alternating steps: 1) Denoising applies a few steps of reverse diffusion to iteratively predict a clean image from the noisy input. 2) Projection enforces data consistency through a proximal update that aligns the output with the degraded observation. By iterating between these steps, DiffPIR achieves effective adversarial defense with strong visual fidelity and minimal architectural overhead.

##### D. Contrastive Learning

Contrastive learning, a self-supervised learning paradigm, can bolster robustness against adversarial attacks by learning discriminative and invariant representations. It trains models to bring augmented views of the same instance closer in the embedding space, while pushing apart unrelated instances, often via a noise contrastive estimation (NCE) loss [28].



Fig. 1: Example of datasets.

Given a batch  $\mathcal{B} = \{x_i\}_{i=1}^N$ , we generate two augmented views per sample  $(\tilde{x}_i, \tilde{x}'_i)$  and optimize the InfoNCE loss:

$$\mathcal{L}_{\text{contrast}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}'_i)/\tau)}{\sum_{k=1}^K \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)} \quad (10)$$

where,  $\mathbf{z}_i = g_\phi(f_\theta(\tilde{x}_i))$  is the projected embedding, with  $f_\theta$  as the encoder (e.g., YOLOv8 backbone) and  $g_\phi$  as the projection head (MLP);  $\text{sim}(\mathbf{u}, \mathbf{v}) = \mathbf{u}^\top \mathbf{v} / (\|\mathbf{u}\| \|\mathbf{v}\|)$  denotes cosine similarity;  $\tau > 0$  is a temperature hyperparameter controlling the softmax sharpness;  $K = |\mathcal{B}| - 1$  is the number of negative pairs per anchor (in-batch negatives);  $\mathbb{1}_{[k \neq i]}$  ensures only true negatives are considered in the denominator. We conducted quantitative evaluations of our contrastive learning-enhanced YOLOv8n model on the Traffic Signs Detection dataset.

## V. EXPERIMENTAL EVALUATION

### A. Dataset

We conduct experimental evaluations on both classification and regression tasks. For the classification task, we use YOLOv8, a widely used model in both academia and industry [7], and stop sign images from the Traffic Signs Detection dataset [29]. The relative distance to the leading vehicle is the key data to ensure driving safety. Therefore, for the regression task, we use the Supercombo model, an end-to-end model used on a production ADS [8], [30], [31], with videos from the Comma2k19 dataset [32] to predict the relative distance to the leading vehicle. An example from each dataset is shown in Fig. 1.

### B. Attack Evaluation

1) *Relative Distance Prediction*: The Supercombo model in OpenPilot is a multitask deep neural network that uses camera images to perform detection of lane lines and lead objects, as well as planning of control outputs, enabling essential functions for autonomous driving systems. We apply various attack methods to generate adversarial patches in the region of the leading vehicle in each video frame. The model’s predicted relative distances under attack are then compared to the predictions on clean images in each frame. We evaluate the average prediction error before and after the attacks across different distance ranges.

As shown in Table I, Gaussian noise exhibits the weakest attack effect, with average errors across all distance ranges remaining minimal. This indicates that the model possesses a degree of robustness to random perturbations. Auto-PGD, on the other hand, produces the largest errors in every distance range, with the highest average error reaching 34.45 meters.

Across all attack methods, the prediction error is notably higher at shorter distances, particularly within 20 meters. This suggests that the model is more vulnerable to adversarial

attacks at close range, likely because perturbations occupy a larger visual area when the target vehicle is nearby, thereby enhancing the attack’s effectiveness.

TABLE I: Avg. errors at different ranges (m) under attack

Attack Method	Range (m)			
	[0, 20]	[20, 40]	[40, 60]	[60, 80]
Gaussian Noise	0.30	0.01	0.03	0.14
FGSM	18.34	4.25	3.92	4.65
Auto-PGD	<b>34.45</b>	<b>8.43</b>	<b>8.11</b>	<b>8.49</b>
CAP-Attack	29.62	6.73	6.42	6.83

2) *Stop Sign Detection*: Fig. 2 presents the performance of the YOLOv8 model in detecting stop signs under different attacks. We evaluate the model using three key metrics: mAP@50, Precision, and Recall. Specifically, mAP@50 quantifies the model’s overall detection accuracy at an Intersection-over-Union (IoU) threshold of 0.5, Precision indicates the proportion of correctly identified objects among all detections, and Recall measures the model’s ability to detect all relevant stop signs.

To streamline the evaluation process, we simplify the detection task by configuring YOLOv8 for single-class detection, focusing exclusively on stop signs. As illustrated in the figure, FGSM and Gaussian Noise attacks lead to the most substantial drop in detection performance, particularly in Recall and mAP@50. Interestingly, although Auto-PGD is generally considered a stronger adversarial attack, its effectiveness appears limited in this context.

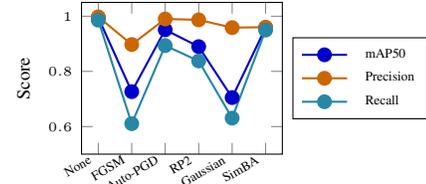


Fig. 2: Performance of stop sign detection with or w/o attacks.

### C. Defense Evaluation

1) *Image Processing*: To evaluate the effectiveness of classical image processing–based defenses against adversarial attacks, we conducted a series of experiments using four commonly applied input-level techniques. The results are summarized in Table II.

In stop sign detection, defense performance is highly dependent on the type of attack. Median Blurring is effective against simple attacks, improving mAP@50 from 70.49% to 94.64% under Gaussian noise and boosting precision from 89.74% to 97.32% under FGSM by smoothing perturbations while preserving structural edges. However, its benefits diminish under stronger attacks such as Auto-PGD, where it yields only marginal gains. In some cases, Randomization and Bit Depth Reduction even degrade performance, likely due to Auto-PGD’s already limited impact in a single-class detection task, which leaves little room for improvement and increases the risk of over-processing. Notably, Randomization

TABLE II: Performance after image processing

Attack	Defense	Avg. Error in Different Range (m)				Stop Sign Detection (%)		
		[0, 20]	[20, 40]	[40, 60]	[60, 80]	mAP50	Prec.	Recall
Gaussian	None	0.30	0.01	0.03	0.14	70.49	95.86	63.06
	Median Blurring	<b>-0.02</b>	-0.12	-0.89	-1.89	94.64	<b>98.47</b>	89.86
	Randomization	4.51	-0.21	-11.44	<b>-24.18</b>	82.95	91.76	74.55
	Bit Depth	1.66	-0.53	-2.33	-2.57	70.17	96.88	62.61
FGSM	None	18.34	4.25	3.92	4.65	72.65	89.74	61.04
	Median Blurring	17.60	6.19	0.66	1.02	85.49	<b>97.32</b>	75.00
	Randomization	<b>4.62</b>	0.40	-10.23	<b>-21.99</b>	75.58	89.91	70.49
	Bit Depth	19.47	6.04	-0.25	-0.89	74.04	87.01	63.38
Auto-PGD	None	34.45	8.43	8.11	8.49	95.09	99.00	89.34
	Median Blurring	25.57	6.65	1.37	2.20	98.84	<b>99.76</b>	94.80
	Randomization	<b>5.04</b>	-0.20	-11.56	<b>-21.25</b>	93.37	94.97	89.26
	Bit Depth	28.87	7.89	0.75	0.94	93.82	98.61	88.06
CAP/ RP <sub>2</sub>	None	29.62	6.73	6.42	6.83	88.97	98.67	83.77
	Median Blurring	25.55	6.16	1.58	1.98	89.54	97.93	83.56
	Randomization	<b>4.73</b>	0.16	-11.44	<b>-21.38</b>	87.39	<b>98.93</b>	81.30
	Bit Depth	28.17	6.92	1.00	0.64	88.97	98.67	83.77

underperforms under Gaussian noise, as its stochastic transformations (e.g., resizing, padding, noise injection) disrupt feature stability in already noisy inputs, reducing detection accuracy.

In regression tasks, image processing defenses also reduce adversarial prediction errors, but their effectiveness varies with object distance. Randomization is most effective at close range, reducing average error from 34.45 meters to 5.04 meters under Auto-PGD. However, beyond 40 meters, its performance drops significantly and even increases error. This is due to the sparsity and low detail of distant vehicle pixels, where random transformations can distort or erase key features, impeding depth estimation. Bit-Depth Reduction and Median Blurring offer moderate improvements, but they are less effective at close range, where stronger perturbations dominate.

These findings suggest that no single preprocessing method provides consistent robustness across different attacks and task conditions. Addressing these challenges may require combining complementary preprocessing techniques or adopting multi-model fusion strategies that account for task- and context-specific vulnerabilities.

2) *Adversarial Training*: We perform adversarial training on the model using different adversarial examples as training sets, and then evaluate the retrained model under various attack scenarios. The results are summarized in Table III.

In order to evaluate transferability, we retrained separate models using adversarial examples generated from each attack (416 stop sign images or 9600 driving video frames), and tested each retrained model on the other three adversarial examples separately. Additionally, to construct a diverse and representative adversarial dataset, we randomly selected 25% of the attacked examples (104 images or 2400 frames) from each of the four attacks to be the mixed training set and randomly selected another 25% of attacked images from each attack to be the test set. Therefore, the mixed adversarial dataset is consistent with each individual dataset and is used to train and test another model for evaluating robustness under compound adversarial conditions.

YOLOv8 models trained with different adversarial examples show varying levels of robustness when evaluated across multiple attack methods. Specifically, the model trained with RP<sub>2</sub> examples performs particularly poorly against FGSM inputs, achieving the lowest scores across all evaluations, with mAP50 dropping to 40.78%, precision to 67.01%, and recall

TABLE III: Performance after adversarial training

Adversarial	Attack	Avg. Error in Different Range (m)				Stop Sign Detection (%)		
		[0, 20]	[20, 40]	[40, 60]	[60, 80]	mAP50	Precision	Recall
Gaussian	FGSM	4.21	3.86	3.11	0.31	89.27	<b>99.76</b>	93.69
	Auto-PGD	7.43	6.83	4.94	7.01	95.62	99.72	98.20
	CAP/RP <sub>2</sub>	5.90	5.69	3.96	0.80	83.98	94.70	85.59
	Mixed	-	-	-	-	83.32	97.69	93.85
	Gaussian	-0.06	-0.08	-0.05	-0.20	94.46	99.32	98.60
FGSM	Auto-PGD	<b>16.83</b>	4.38	2.77	3.97	95.13	<b>99.54</b>	98.39
	CAP/RP <sub>2</sub>	12.51	3.50	2.24	2.08	87.83	97.11	87.61
	Mixed	-	-	-	-	86.68	98.29	95.78
	Gaussian	-0.09	0.01	0.10	0.05	98.85	<b>99.72</b>	98.65
	FGSM	2.92	1.73	1.76	1.87	85.60	98.57	87.61
Auto-PGD	CAP/RP <sub>2</sub>	6.50	2.99	2.34	2.50	85.61	97.65	84.40
	Mixed	-	-	-	-	82.70	97.88	91.94
	Gaussian	0.18	0.53	-0.02	-1.43	55.13	92.20	54.50
	FGSM	3.31	3.04	1.91	-1.43	40.78	67.01	60.40
	Auto-PGD	7.78	5.24	3.19	5.10	88.13	<b>98.23</b>	87.56
CAP/RP <sub>2</sub>	Mixed	-	-	-	-	67.68	92.44	74.69
	Gaussian	0.25	0.02	-0.02	<b>-43.04</b>	90.57	98.62	96.70
	FGSM	4.17	1.61	1.41	0.16	88.63	98.14	95.19
	Auto-PGD	5.84	2.91	3.42	1.47	87.90	99.09	97.64
	CAP/RP <sub>2</sub>	4.89	2.39	2.86	0.79	81.95	97.88	88.96
Mixed	Mixed	-	-	-	-	88.47	<b>99.42</b>	94.92

to 60.40%. This sharp decline highlights a key vulnerability: RP<sub>2</sub>-trained models fail to generalize to gradient-based attacks, indicating overfitting to the RP<sub>2</sub> attack pattern and poor cross-attack robustness. In contrast, models trained with gradient-based perturbations like FGSM or Auto-PGD exhibit more consistent and balanced performance across diverse adversarial inputs.

In relative distance prediction, we retrain the model in the same way, and adversarial training proves effective, particularly at close range. When trained with a mix of adversarial examples, the model’s average error under Auto-PGD within 0–20 meters drops sharply from 34.45 meters to just 5.84 meters. Compared to single-attack training, mixed adversarial training offers more balanced robustness across diverse attack types, highlighting its potential as a general-purpose defense strategy. However, this improvement comes at a cost. The same model exhibits significantly larger errors at longer distances, with a maximum mean error reaching -43.04 meters. This behavior suggests that the adversarial loss reshapes the feature space, optimizing robustness in the near field (0–20m) but limiting generalization beyond that range. To address this, future defenses should explore distance-aware sampling, loss weighting, or multi-scale adversarial training.

3) *Contrastive Learning*: We apply contrastive learning to enhance the feature representation of a YOLOv8 model, aiming to improve robustness in stop sign detection. This self-supervised learning approach utilizes a projection head with batch normalization and dropout, and employs a multi-positive contrastive loss with a margin to promote better feature separation. The training and test sets are the same as those for adversarial training. Performance results across different attacks are summarized in Table IV.

While contrastive learning yields some improvements, the gains are modest. This is likely due to contrastive learning’s emphasis on feature invariance, which reduces sensitivity to input variations but does not explicitly target adversarial robustness [33]. In contrast, adversarial training with a single attack type often leads to overfitting, improving robustness against that attack while reducing generalization to others. Additionally, the contrastive-trained model maintains strong

TABLE IV: Performance after contrastive learning

Adv. Example	Attack Method	mAP50 (%)	Precision (%)	Recall (%)
Gaussian Noise	Clean	99.49	<b>100.00</b>	97.53
	FGSM	78.11	68.93	59.16
	Auto-PGD	96.24	96.01	95.17
	RP <sub>2</sub>	92.22	98.92	79.31
	SimBA	99.41	99.53	99.05
FGSM	Clean	99.49	99.06	98.02
	Gaussian Noise	80.36	91.53	66.26
	Auto-PGD	98.50	98.14	97.15
	RP <sub>2</sub>	93.29	98.26	<b>100.00</b>
	SimBA	99.41	99.84	98.37
Auto-PGD	Clean	99.49	99.06	98.17
	Gaussian Noise	79.82	98.44	60.47
	FGSM	79.67	95.12	56.60
	RP <sub>2</sub>	90.22	99.86	81.32
	SimBA	99.10	<b>99.90</b>	98.51
RP <sub>2</sub>	Clean	99.39	99.76	97.01
	Gaussian Noise	79.01	97.69	59.34
	FGSM	75.74	95.79	53.77
	Auto-PGD	91.02	93.49	84.01
	SimBA	92.11	<b>100.00</b>	83.51
SimBA	Clean	99.39	<b>99.61</b>	98.11
	Gaussian Noise	79.58	75.34	53.83
	FGSM	75.98	94.30	52.11
	Auto-PGD	92.50	82.41	81.67
	RP <sub>2</sub>	92.41	98.82	84.90

and consistent performance under FGSM, suggesting that feature-level invariance may offer some resilience to simpler, gradient-based perturbations.

4) *Diffusion Model*: Table V presents the performance of the diffusion model in mitigating adversarial perturbations. The results demonstrate that diffusion-based image reconstruction serves as an effective defense strategy, particularly in relative distance prediction tasks. For instance, under Auto-PGD, the average prediction error within a 20-meter range drops dramatically from 34.45 meters to just 4.98 meters. However, in cases like Gaussian noise, where the original attack is relatively weak, the diffusion process can inadvertently degrade performance, introducing unnecessary changes that lead to increased prediction errors. Similarly, at longer distances, diffusion-repaired images often result in negatively biased predictions, with average errors becoming negative. This suggests that the model systematically underestimates distances after reconstruction, possibly due to subtle distortions or overcorrections introduced by the generative process.

In stop sign detection, the YOLOv8 model achieves a precision of over 99% across all attacks, and reaches 100% precision under FGSM, showing the strength of diffusion models in restoring adversarially perturbed inputs for classification.

TABLE V: Performance after diffusion model cleaning

Attack Method	Avg. Error in Different Range (m)				Stop Sign Detection (%)		
	[0, 20]	[20, 40]	[40, 60]	[60, 80]	mAP50	Precision	Recall
Gaussian	-1.06	-1.33	-4.32	-5.27	99.45	99.50	97.99
FGSM	6.60	-1.06	-3.31	-5.03	97.81	<b>100.00</b>	93.98
Auto-PGD	4.98	-1.05	-3.90	-3.77	99.50	99.96	99.26
CAP/RP <sub>2</sub>	<b>8.42</b>	-0.67	-2.65	-4.08	93.74	99.33	89.95
SimBA	-	-	-	-	99.50	99.97	99.26

## VI. DISCUSSION

The evaluation of defense methods indicates that no single defense approach can fully protect against all types of attacks. For instance, the Median Blurring enhances the model’s robustness to a certain extent against multiple attacks, such as

mitigating errors caused by FGSM and Auto-PGD attacks by more than 10 meters, but is less effective against more subtle or dynamic threats, mitigating only about 4 meters under CAP-Attack. Conversely, methods such as diffusion models and contrastive learning show promising defense potential, yet they still require further optimization to minimize their negative impact on model accuracy in the absence of attacks.

Additionally, time overhead is another factor that must be considered. In adversarial training and contrastive learning, the model has been retrained, and they do not require additional time for processing. While in image processing methods, on average, it takes about 20ms to process each frame or image. Most notably, DiffPIR takes 1-2 seconds to repair an attacked image, which is an unacceptable overhead. Because DiffPIR wasn’t designed with the real-time scenario and is bloated, optimizing it for real-time applications deserves further study.

Not only from the experimental results, but also from the time cost, we can observe that adversarial training is a straightforward and effective defense method. In particular, training with mixed adversarial samples achieves a better balance across different attacks but may lead to over-defense, which reduces accuracy in long-distance prediction tasks. Future research should focus on refining adversarial training strategies, particularly in terms of selecting and designing adversarial samples to strike an optimal balance between generalization performance and defensive effectiveness.

## VII. RELATED WORK

Many studies have attempted to attack various perception components of the autonomous driving system, such as LiDAR [34], [35], traffic sign recognition [36]–[39], road lane detection [40], [41], trajectory prediction [42], [43], or vehicle detection [44]–[46]. To evaluate and improve model robustness [47], numerous benchmarks have been proposed [48], [49] as well as many defense methods, such as high-level representation guided denoiser [50], convolutional sparse coding [51], perturbation rectifying network [52], and runtime safety monitoring [53], [54] and interventions [55]. However, most of these strategies [56] either do not directly study the perception module or are still tailored to image classification and evaluated with an offline image dataset. In this work, we revisit the adversarial robustness of deep learning models in the context of autonomous driving for both object detection and regression tasks from attack and defense perspectives.

## VIII. CONCLUSION

In this study, we systematically investigate the vulnerability of ADS perception models to various adversarial attacks and evaluate multiple defense strategies. Our experimental results demonstrate that adversarial attacks significantly affect the model’s classification task and regression tasks, particularly at close range. Among the defense methods evaluated, different approaches show effectiveness under specific conditions, underscoring the urgent need to enhance the robustness of ADS perception models through innovative and adaptive defense strategies.

## REFERENCES

- [1] J. Zhao, W. Zhao, B. Deng, Z. Wang, F. Zhang, W. Zheng, W. Cao, J. Nan, Y. Lian, and A. F. Burke, "Autonomous driving system: A comprehensive survey," *Expert Systems with Applications*, vol. 242, p. 122836, 2024.
- [2] Tesla, Inc. (2024) Model y owner's manual. Accessed: 2025-04-07. [Online]. Available: [https://www.tesla.com/ownersmanual/modely/en\\_us/GUID-2CB60804-9CEA-4F4B-8B04-09B991368DC5.html](https://www.tesla.com/ownersmanual/modely/en_us/GUID-2CB60804-9CEA-4F4B-8B04-09B991368DC5.html)
- [3] L.-H. Wen and K.-H. Jo, "Deep learning-based perception systems for autonomous driving: A comprehensive survey," *Neurocomputing*, vol. 489, pp. 255–270, 2022.
- [4] Y. Cao, C. Xiao, B. Cyr, Y. Zhou, W. Park, S. Rampazzi, Q. A. Chen, K. Fu, and Z. M. Mao, "Adversarial sensor attack on lidar-based perception in autonomous driving," in *Proceedings of the 2019 ACM SIGSAC conference on computer and communications security*, 2019, pp. 2267–2281.
- [5] L. Chi, M. Msahli, Q. Zhang, H. Qiu, T. Zhang, G. Memmi, and M. Qiu, "Adversarial attacks on autonomous driving systems in the physical world: a survey," *IEEE Transactions on Intelligent Vehicles*, 2024.
- [6] H. Wu, S. Yunas, S. Rowlands, W. Ruan, and J. Wahlström, "Adversarial driving: Attacking end-to-end autonomous driving," in *2023 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2023, pp. 1–7.
- [7] Ultralytics, "Yolov8," 2025, accessed: 2025-05-06. [Online]. Available: <https://docs.ultralytics.com/models/yolov8/>
- [8] Comma.ai, "OpenPilot." [Online]. Available: <https://github.com/commaai/openpilot>
- [9] Wikipedia contributors. (2024) Gaussian noise - wikipedia. [Online]. Available: [https://en.wikipedia.org/wiki/Gaussian\\_noise](https://en.wikipedia.org/wiki/Gaussian_noise)
- [10] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," 2015. [Online]. Available: <https://arxiv.org/abs/1412.6572>
- [11] F. Croce and M. Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks," 2020. [Online]. Available: <https://arxiv.org/abs/2003.01690>
- [12] C. Guo, J. R. Gardner, Y. You, A. G. Wilson, and K. Q. Weinberger, "Simple black-box adversarial attacks," 2019. [Online]. Available: <https://arxiv.org/abs/1905.07121>
- [13] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1625–1634.
- [14] X. Zhou, A. Chen, M. Kouzel, H. Ren, M. McCarty, C. Nita-Rotaru, and H. Alemzadeh, "Runtime Stealthy Perception Attacks against DNN-Based Adaptive Cruise Control Systems," in *ACM Asia Conference on Computer and Communications Security (ASIA CCS)*, 2025.
- [15] D. J. Yeong, G. Velasco-Hernandez, J. Barry, and J. Walsh, "Sensor and sensor fusion technology in autonomous vehicles: A review," *Sensors*, vol. 21, no. 6, p. 2140, 2021.
- [16] E. Ward and J. Folkesson, "Vehicle localization with low cost radar sensors," in *2016 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2016, pp. 864–870.
- [17] A. Ruta, F. Porikli, S. Watanabe, and Y. Li, "In-vehicle camera traffic sign detection and recognition," *Machine Vision and Applications*, vol. 22, pp. 359–375, 2011.
- [18] K. Behrendt, L. Novak, and R. Botros, "A deep learning approach to traffic lights: Detection, tracking, and classification," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 1370–1377.
- [19] B. Coifman, D. Beymer, P. McLauchlan, and J. Malik, "A real-time computer vision system for vehicle tracking and traffic surveillance," *Transportation Research Part C: Emerging Technologies*, vol. 6, no. 4, pp. 271–288, 1998.
- [20] J. Zhou, "A review of lidar sensor technologies for perception in automated driving," *Academic Journal of Science and Technology*, vol. 3, no. 3, pp. 255–261, 2022.
- [21] R. R. Wiyatno, A. Xu, O. Dia, and A. De Berker, "Adversarial examples in modern machine learning: A review," *arXiv preprint arXiv:1911.05268*, 2019.
- [22] B. Badjie, J. Cecilio, and A. Casimiro, "Adversarial attacks and countermeasures on image classification-based deep learning models in autonomous driving systems: A systematic review," *ACM Computing Surveys*, vol. 57, no. 1, pp. 1–52, 2024.
- [23] V. Zantedeschi, M.-I. Nicolae, and A. Rawat, "Efficient defenses against adversarial attacks," in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017, pp. 39–49.
- [24] W. Xu, D. Evans, and Y. Qi, "Feature squeezing: Detecting adversarial examples in deep neural networks," in *Proceedings 2018 Network and Distributed System Security Symposium*, ser. NDSS 2018. Internet Society, 2018. [Online]. Available: <http://dx.doi.org/10.14722/ndss.2018.23198>
- [25] C. Xie, J. Wang, Z. Zhang, Z. Ren, and A. Yuille, "Mitigating adversarial effects through randomization," 2018. [Online]. Available: <https://arxiv.org/abs/1711.01991>
- [26] M. Zhao, L. Zhang, J. Ye, H. Lu, B. Yin, and X. Wang, "Adversarial training: A survey," 2024. [Online]. Available: <https://arxiv.org/abs/2410.15042>
- [27] Y. Zhu, K. Zhang, J. Liang, J. Cao, B. Wen, R. Timofte, and L. V. Gool, "Denosing diffusion models for plug-and-play image restoration," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops (NTIRE)*, 2023.
- [28] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," 2020. [Online]. Available: <https://arxiv.org/abs/2002.05709>
- [29] P. K. Darabi, "Car detection dataset," <https://www.kaggle.com/datasets/pkdarabi/cardetection>, 2022, accessed: 2025-04-06.
- [30] A. Schmedding, P. Schowitz, X. Zhou, Y. Lu, L. Yang, H. Alemzadeh, and E. Smirni, "Strategic resilience evaluation of neural networks within autonomous vehicle software," in *International Conference on Computer Safety, Reliability, and Security*. Springer, 2024, pp. 33–48.
- [31] X. Zhou, A. Schmedding, H. Ren, L. Yang, P. Schowitz, E. Smirni, and H. Alemzadeh, "Strategic Safety-Critical Attacks against an Advanced Driver Assistance System," in *2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 2022, pp. 79–87.
- [32] H. Schafer, E. Santana, A. Haden, and R. Biasini, "A commute in data: The comma2k19 dataset," 2018.
- [33] A. Foster, R. Pukdee, and T. Rainforth, "Improving transformation invariance in contrastive representation learning," 2021. [Online]. Available: <https://arxiv.org/abs/2010.09515>
- [34] J. Chen, M. Su, S. Shen, H. Xiong, and H. Zheng, "Poba-ga: Perturbation optimized black-box adversarial attacks via genetic algorithm," *Computers & Security*, vol. 85, pp. 89–106, 2019.
- [35] K. Yang, T. Tsai, H. Yu, M. Panoff, T.-Y. Ho, and Y. Jin, "Robust roadside physical adversarial attack against deep learning in lidar perception modules," in *Proceedings of the 2021 ACM Asia Conference on Computer and Communications Security*, 2021, pp. 349–362.
- [36] X. Wei, Y. Guo, J. Yu, and B. Zhang, "Simultaneously optimizing perturbations and positions for black-box adversarial patch attacks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 7, pp. 9041–9054, 2022.
- [37] X. Wei, Y. Guo, and J. Yu, "Adversarial sticker: A stealthy attack method in the physical world," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 2711–2725, 2022.
- [38] C. Hu and W. Shi, "Adversarial color film: Effective physical-world attack to dnn," *arXiv preprint arXiv:2209.02430*, 2022.
- [39] X. Yang, W. Liu, S. Zhang, W. Liu, and D. Tao, "Targeted attention attack on deep learning models in road sign recognition," *IEEE Internet of Things Journal*, vol. 8, no. 6, pp. 4980–4990, 2020.
- [40] P. Jing, Q. Tang, Y. Du, L. Xue, X. Luo, T. Wang, S. Nie, and S. Wu, "Too good to be safe: Tricking lane detection in autonomous driving with crafted perturbations," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 3237–3254.
- [41] A. Boloor, K. Garimella, X. He, C. Gill, Y. Vorobeychik, and X. Zhang, "Attacking vision-based perception in end-to-end autonomous driving models," *Journal of Systems Architecture*, vol. 110, p. 101766, 2020.
- [42] R. Muller, Y. Man, Z. B. Celik, M. Li, and R. Gerdes, "Physical hijacking attacks against object trackers," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, 2022, pp. 2309–2322.
- [43] Y. Cao, C. Xiao, A. Anandkumar, D. Xu, and M. Pavone, "Advdo: Realistic adversarial attacks for trajectory prediction," in *European Conference on Computer Vision*. Springer, 2022, pp. 36–52.
- [44] J. Tu, M. Ren, S. Manivasagam, M. Liang, B. Yang, R. Du, F. Cheng, and R. Urtasun, "Physically realizable adversarial examples for lidar object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 13 716–13 725.

- [45] K. Yang, T. Tsai, H. Yu, T.-Y. Ho, and Y. Jin, "Beyond digital domain: Fooling deep learning based recognition system in physical world," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 01, 2020, pp. 1088–1095.
- [46] J. Wang, A. Liu, Z. Yin, S. Liu, S. Tang, and X. Liu, "Dual attention suppression attack: Generate adversarial camouflage in physical world," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 8565–8574.
- [47] X. Zhou, M. Kouzel, and H. Alemzadeh, "Robustness testing of data and knowledge driven anomaly detection in cyber-physical systems," in *2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*. IEEE, 2022, pp. 44–51.
- [48] A. E. Cinà, J. Rony, M. Pintor, L. Demetrio, A. Demontis, B. Biggio, I. B. Ayed, and F. Roli, "Attackbench: Evaluating gradient-based attacks for adversarial examples," 2024. [Online]. Available: <https://arxiv.org/abs/2404.19460>
- [49] F. Croce and M. Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks," in *ICML*, 2020.
- [50] F. Liao, M. Liang, Y. Dong, T. Pang, X. Hu, and J. Zhu, "Defense against adversarial attacks using high-level representation guided denoiser," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [51] B. Sun, N.-H. Tsai, F. Liu, R. Yu, and H. Su, "Adversarial defense by stratified convolutional sparse coding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [52] N. Akhtar, J. Liu, and A. Mian, "Defense against universal adversarial perturbations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [53] X. Zhou, B. Ahmed, J. H. Aylor, P. Asare, and H. Alemzadeh, "Hybrid knowledge and data driven synthesis of runtime monitors for cyber-physical systems," *IEEE Transactions on Dependable and Secure Computing*, vol. 21, no. 1, pp. 12–30, 2023.
- [54] —, "Data-driven design of context-aware monitors for hazard prediction in artificial pancreas systems," in *2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 2021, pp. 484–496.
- [55] C. Chen, G. Xiao, D. Lee, L. Yang, E. Smirni, H. Alemzadeh, and X. Zhou, "Safety interventions against adversarial patches in an open-source driver assistance system," *to appear in the 55th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, 2025.
- [56] A. D. M. Ibrahim *et al.*, "Adversarial attacks and defenses in deep learning for autonomous vehicles: A systematic review from a safety perspective," *Artificial Intelligence Review*, 2024.