

On the Natural Robustness of Vision-Language Models Against Visual Perception Attacks in Autonomous Driving

Pedram MohajerAnsari¹, Amir Salarpour¹, Michael Kühr², Siyu Huang¹, Mohammad Hamad², Sebastian Steinhorst², Habeeb Olufowobi³, Mert D. Pesé¹

¹Clemson University, Clemson, SC, USA

²Technical University of Munich, Munich, Germany

³University of Texas at Arlington, Arlington, TX, USA

{pmohaje, asalarp, siyuh, mpese}@clemson.edu
 {michael.kuehr, mohammad.hamad, sebastian.steinhorst}@tum.de
 habeeb.olufowobi@uta.edu

Abstract

Autonomous vehicles (AVs) rely on deep neural networks (DNNs) for critical tasks such as traffic sign recognition (TSR), automated lane centering (ALC), and vehicle detection (VD). However, these models are vulnerable to attacks that can cause misclassifications and compromise safety. Traditional defense mechanisms, including adversarial training, often degrade benign accuracy and fail to generalize against unseen attacks. In this work, we introduce Vehicle Vision Language Models (V^2LMs), fine-tuned Vision-Language Models (VLMs) specialized for AV perception. Our findings demonstrate that V^2LMs inherently exhibit superior robustness against unseen attacks without requiring adversarial training, maintaining significantly higher accuracy than conventional DNNs under adversarial conditions. We evaluate two deployment strategies: Solo Mode, where individual V^2LMs handle specific perception tasks, and Tandem Mode, where a single unified V^2LM is fine-tuned for multiple tasks simultaneously. Experimental results reveal that DNNs suffer performance drops of 33%–46% under attacks, whereas V^2LMs maintain adversarial accuracy with reductions of less than 8% on average. The Tandem Mode further offers a memory-efficient alternative while achieving comparable robustness to the Solo Mode. Also, we explore the integration of V^2LMs as parallel components to AV perception to enhance resilience against adversarial threats. Our results suggest that V^2LMs offer a promising path toward more secure and resilient AV perception systems.¹

¹The implementation code of this work is available at <https://github.com/pedram-mohajer/V2LM>

■ Victim Path without V^2LM ■ NPC Path ■ Victim Path with V^2LM

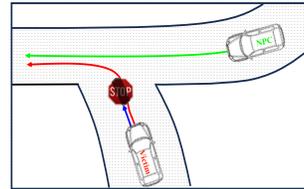
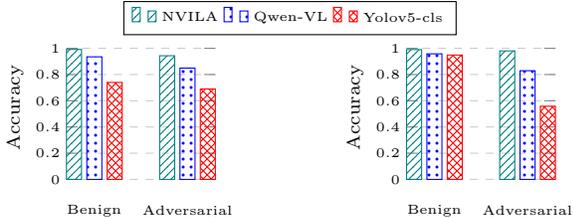


Figure 1: An attacker casts a shadow on the stop sign, causing the AV with a traditional TSR algorithm to misclassify it and continue driving (red path), leading to a collision. In contrast, with V^2LM , the AV remains robust and correctly stops (blue path) despite the attack.

1 Introduction

Autonomous vehicles (AVs) are built on automated driving systems (ADS), consisting of modular subsystems including perception, planning, and control. The perception module—significantly enhanced by advancements in deep neural network (DNN) models for tasks such as image classification and object detection [59]—reads and processes data from sensors such as cameras and LiDAR; this information is then sent to the planning module, which makes decisions on navigation and maneuvers, and subsequently to the control module, executing these decisions to control the vehicle’s movement [49].

While these DNN models effectively recognize, monitor, and predict the movements of nearby objects [52, 67], they are also vulnerable to physical attacks that manipulate real-world environments to deceive perception systems [69, 24, 38, 39]. Such physical attacks pose a significant threat to AV perception by causing misclassification and unsafe driving



(a) **After Adversarial Training.** Adversarial training improves robustness but reduces benign accuracy, especially for YOLOv5-cls. (b) **Before Adversarial Training.** Fine-tuned VLMs (NVILA, Qwen-VL) exhibit strong robustness even without adversarial training, while traditional DNNs like YOLOv5-cls show greater vulnerability and poorer generalization.

Figure 2: Accuracy comparison of NVILA, Qwen-VL, and YOLOv5-cls on benign and adversarial inputs before and after adversarial training. VLMs maintain high robustness with minimal trade-off, while traditional DNNs like YOLOv5-cls show greater vulnerability and poorer generalization.

behaviors, as demonstrated by methods such as Drp-Attack [48], ControlLoc [35], and SlowTrack [34].

Various defense methods have been proposed to mitigate such attacks: defensive distillation reduces model sensitivity but struggles to generalize across diverse perturbations [41]; input transformations can suppress adversarial noise but often degrade clean data quality, lowering accuracy [20]; and provable defenses, though theoretically robust, are computationally expensive and difficult to scale [60]. Among these, adversarial training is the most widely adopted defense. Although it improves robustness against specific attacks, it often reduces generalization and standard accuracy on clean data [43], as shown in Figure 2a, which is especially problematic in AV contexts where high performance under normal conditions is essential for safe and reliable operation [58].

To address these limitations and achieve both robust and generalizable models for AV perception tasks, this work proposes a novel finding: *fine-tuned VLMs inherently exhibit robustness against unseen adversarial attacks on AV perception systems even without adversarial training, significantly outperforming DNNs*. As shown in Figure 2(b), VLMs such as Qwen-VL [4] and NVILA [32] demonstrate strong resilience, achieving 82.93% and 98.04% accuracy respectively against adversarial examples without adversarial training. In contrast, the YOLOv5-cls [26, 44] model suffers a drastic drop from 94.88% accuracy on benign inputs to just 55.82% on adversarial ones. Based on this observation, we introduce fine-tuned VLMs for the first time as a robust AV perception module, naming them Vehicle Vision Language Models (V^2LMs) and comprehensively evalu-

ating them in the AV context.

This paper evaluates six VLMs, namely LLaVA-7B, LLaVA-13B-LoRA [30], MoE-LLaVA [28], MobileVLM [15], Qwen-VL-7B, and NVILA-8B for their potential as auxiliary in-vehicle components to enhance perception performance in the presence of attacks targeting perception algorithms. The first four models are based on LLaMA [56], while the last is based on Qwen [3], selected for its state-of-the-art performance in multimodal tasks and its efficiency in offline operation—an essential feature for our application. In addition to these VLMs, the study also evaluates task-specific DNN models, including YOLOv5-cls [26, 44] for Traffic Sign Recognition (TSR), CLRRerNet [21] for Automated Lane Centering (ALC), and YOLOv5-dt for Vehicle Detection (VD), as strong baselines under both benign and adversarial conditions.

The study begins by evaluating the VLMs’ zero-shot performance on key AV tasks, including (i) TSR, (ii) ALC, and (iii) VD as part of Object Detection (OD). Subsequently, they are fine-tuned specifically for these tasks to better align them with AV applications, with their performance reassessed to determine improvements (RQ1). We refer to each fine-tuned model as a Vehicle Vision-Language Model (V^2LM), highlighting its specialization for AV tasks. Additionally, the task-specific DNN-based models are fine-tuned on the same training dataset as the VLMs, ensuring a fair comparison of their capabilities under identical conditions.

Then, the effectiveness and resilience of both DNN models and V^2LMs are tested using *unseen* adversarial examples (AEs) against three distinct attacks targeting AV perception algorithms (RQ2): (1) Robust and Accurate UV-map-based Camouflage attack (RAUCA) to deceive VD algorithms [70], (2) a physical-world adversarial attack known as the Dirty Road Patch (DRP-Attack) to compromise DNN-based automated lane centering (ALC) models [48], and (3) shadows cast on traffic signs to attack TSR algorithms [69]. As shown in Figure 1, the V^2LM improves robustness against an adversarial attack on a traffic sign, helping the AV avoid a collision.

The study then compares two distinct designs for utilizing V^2LMs and evaluates their performance on the three aforementioned AV tasks and AEs (RQ3). The first design, termed *Solo Mode*, involves separate V^2LMs , each fine-tuned individually for one of the AV tasks. The second design, named *Tandem Mode*, uses a single V^2LM fine-tuned simultaneously for all three AV tasks, aiming to improve robustness against various attack types. By comparing these two designs, we aim to determine whether a

unified model can achieve robustness across multiple tasks and adversarial conditions as effectively as task-specific models.

Finally, we discuss the feasibility of integrating a V^2LM as a parallel component within AV systems to bolster the perception module against attacks. Given the critical need for low latency to ensure real-time decision-making in autonomous driving, V^2LM 's latency is evaluated relative to the perception system and human reaction times. This assessment aims to determine whether V^2LM can operate within strict real-time constraints without adversely impacting the vehicle's responsiveness.

This paper makes the following contributions:

- This work proposes a novel finding that fine-tuned Vision Language Models (VLMs) inherently exhibit superior robustness against *unseen* adversarial attacks compared to task-specific DNNs. We highlight a critical limitation of traditional defense methods, such as adversarial training: while adversarial training aims to improve robustness, it significantly degrades the benign accuracy of traditional DNNs while providing only limited improvements against adversarial examples. In contrast, VLM-based models achieve strong adversarial robustness while maintaining high benign accuracy.
- We introduce Vehicle Vision Language Models (V^2LMs), fine-tuned VLMs specifically for AV perception tasks: TSR, ALC, and VD. We further propose two deployment strategies: *Solo Mode*, where separate VLM are fine-tuned individually for each task, and *Tandem Mode*, employing a single unified V^2LM across multiple perception tasks in same time.
- We conduct comprehensive experiments to evaluate the robustness of traditional DNN models and V^2LMs under adversarial conditions. Traditional DNN models experience performance drops of 33%–46% under attacks, whereas V^2LMs achieve significantly smaller reductions (less than 8% on average), consistently maintaining high adversarial accuracy without relying on additional defense mechanisms.

2 Related Work

Large Language Models in Autonomous Driving. Recent work has demonstrated the potential of LLMs in the context of AVs, particularly enhancing perception, control, and motion planning tasks.

Table 1: Zero-shot performance of VLMs on AV tasks (Solo Design).

Task	Model	Accuracy	F1-Score	Precision	Recall
TSR	LLaVA-7B	8.23%	9.19%	2.38%	2.19%
	LLaVA-13B-LoRA	13.51%	15.03%	4.93%	6.04%
	MoE-LLaVA	1.19%	1.63%	2.06%	0.48%
	MobileVLM	1.27%	0.95%	1.15%	0.81%
	Qwen-7B	7.01%	1.73%	4.52%	6.08%
	NVILA-8B	31.24%	29.80%	29.24%	30.31%
ALC	LLaVA-7B	37.16%	36.56%	42.71%	29.33%
	LLaVA-13B-LoRA	38.20%	38.12%	45.86%	25.64%
	MoE-LLaVA	35.41%	28.84%	29.68%	21.39%
	MobileVLM	20.71%	24.60%	28.38%	21.71%
	Qwen-7B	26.96%	11.67%	7.45%	26.97%
	NVILA-8B	50.91%	51.18%	51.44%	50.91%
VD	LLaVA-7B	77.91%	64.01%	88.09%	51.61%
	LLaVA-13B-LoRA	91.26%	84.46%	89.82%	80.23%
	MoE-LLaVA	75.61%	65.16%	87.95%	50.46%
	MobileVLM	60.64%	67.83%	71.43%	64.57%
	Qwen-7B	75.44%	82.57%	94.38%	75.44%
	NVILA-8B	92.06%	92.01%	91.82%	92.06%

Regarding perception systems, LLMs utilize external APIs to access real-time information sources, including traffic reports, and weather updates, which significantly enrich the vehicle's ability to gain a comprehensive understanding of its environment [16]. Aldeen *et al.* [2] investigate the application of Large Multimodal Models (LMMs) for enhancing the cybersecurity of AVs. Concerning control, LLMs enable the adjustment of control settings according to driver preferences, thereby personalizing the driving experience [50].

Additionally, LLMs enhance transparency by providing detailed explanations of each step in the control process. These capabilities extend to navigation improvements; e.g., LLMs effectively process real-time traffic data to identify congested routes and suggest alternative paths, significantly optimizing navigation for both efficiency and safety [53]. For motion planning tasks, LLMs leverage their advanced natural language understanding and reasoning capabilities [37]. This allows for enhanced user-centric communication, enabling passengers to articulate their intentions and preferences in everyday language. Furthermore, LLMs analyze textual data from sources such as maps, traffic reports, and real-time updates to make high-level decisions that optimize route planning [40].

Adversarial Threats and Defenses. Adversarial Examples (AEs) were first defined by Szegedy *et al.* [55] as subtly modified inputs designed to fool DNNs. These minor, often imperceptible alterations can drastically alter a DNN's predictions [51, 9]. Several studies have demonstrated the vulnerability of AV perception systems to AEs. Kong *et al.* [27] introduced PhysGAN, a GAN-based framework that generates AEs resilient in the physical world, capable of misleading ADS throughout an entire trajectory. Similarly, Eykholt *et al.* [18] showed that even plac-

ing a sticker on a stop sign could deceive the AV’s TSR.

To counter these vulnerabilities, various defense mechanisms have been proposed. Adversarial training, proposed by Madry *et al.* [36], involves training models on both clean and AEs to enhance robustness. Defensive distillation, introduced by Papernot *et al.* [41], aims to reduce the model’s sensitivity to perturbations by smoothing decision boundaries. Robust optimization techniques, discussed by Wong and Kolter [60], improve the model’s resistance to adversarial attacks through advanced optimization methods. Also, input preprocessing techniques, such as those outlined by Guo *et al.* [20], involve denoising or transforming inputs to mitigate the effects of adversarial perturbations.

3 V^2LM as a Defense Mechanism

3.1 Overview

The perception system in AVs is responsible for interpreting the environment through sensor data, such as images, to enable safe and efficient operation. This system performs essential tasks including TSR, ALC, and VD which is part of object detection [65]. TSR enables the vehicle to follow road rules by recognizing traffic signs [57], ALC keeps the vehicle centered within its lane by identifying road markings [17], and VD detects and classifies other vehicles [8, 11]. The perception system (PS) can be represented as a function, shown in Equation 1, which takes an input image $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$ and outputs the results of TSR, ALC, and VD. Specifically, TSR outputs the detected traffic sign’s class and bounding box if a sign is present; ALC provides a classification for the appropriate steering command; and VD returns the bounding box and class of any detected car.

$$PS(I) = \begin{cases} TSR(\mathcal{I}) : Cls_{traffic-sign} \\ ALC(\mathcal{I}) : Cls_{steering} \\ VD(\mathcal{I}) : (Cls_{vehicle}, BBox_{vehicle}) \end{cases} \quad (1)$$

where Cls denotes class and $BBox$ denotes bounding box.

AEs can interfere with the perception system by causing errors in these modules. An AE, \mathcal{I}_{adv} , is crafted by adding a small perturbation δ to an original input \mathcal{I} such that the target module’s output changes undesirably:

$$\mathcal{I}_{adv} = \mathcal{I} + \delta, \quad \text{where } \|\delta\| \leq \epsilon \quad \text{and} \quad f(\mathcal{I}_{adv}) \neq f(\mathcal{I}) \quad (2)$$

where f denotes the specific target module in the AV’s perception system, which can be the TSR, ALC, or VD module. This perturbation δ , constrained by ϵ , represents a small, controlled modification designed to be imperceptible, thereby ensuring that \mathcal{I}_{adv} visually resembles \mathcal{I} . Despite this, the adversarial perturbation may lead to misclassification of traffic signs, incorrect lane-following commands, or faulty detection of objects.

3.2 Vision Language Model (VLM)

VLMs, by integrating visual and textual data [66], have the potential to enhance perception tasks by enabling a more comprehensive understanding of the environment. A pre-trained VLM takes an image and a text prompt as inputs to generate a relevant response:

$$VLM_{pre}(\mathcal{I}, Prompt) \rightarrow Generated-Response \quad (3)$$

This potential arises from their several key strengths: their multimodal learning capability, which allows them to correlate visual and textual information simultaneously [33]; their robustness to variability, which enables them to generalize well across different environments due to extensive training on diverse datasets [46]; their contextual understanding, which leverages textual data to enhance the interpretation of visual scenes [68]; and their comprehensive feature extraction, which combines features from both visual and textual data [5]. For instance, VLMs could enhance ALC by interpreting road markings and reading associated signs. In VD, VLMs could recognize and classify objects like vehicles by analyzing visual data along with bounding box coordinates. Similarly, VLMs could improve TSR by understanding text on signs, such as speed limits or warnings, to ensure the car follows road rules accurately.

3.3 Vehicle Vision Language Model (V^2LM)

Fine-tuning VLMs on AV-specific tasks is essential for optimizing their performance and enabling adaptation to domain-specific challenges such as variations in lighting, road conditions, and traffic scenarios; the resulting fine-tuned VLM is referred to as V^2LM . For LLaVA, this fine-tuning process involves adjustments

to key components. The vision encoder, implemented as CLIP ViT-L/14 [42], extracts visual features from input images. The language model, based on Vicuna (a fine-tuned variant of LLaMA) [14], interprets textual prompts. A cross-modal attention module integrates the modalities, and all components are jointly optimized via visual instruction tuning [31] to align with autonomous vehicle (AV) tasks.

Qwen-VL adopts a modular architecture comprising an OpenCLIP ViT-bigG visual encoder [13], a single-layer cross-attention adapter with learnable queries, and a Qwen-7B language model [3]. The adapter compresses image features to fixed-length sequences, enabling efficient visual grounding and fine-grained perception. Its three-stage training pipeline includes weakly supervised pretraining on large-scale image-text pairs, multi-task visual-language pretraining, and supervised instruction tuning. **NVILA** builds on this by incorporating a “scale-then-compress” design, which first increases spatial and temporal input resolution and then compresses visual tokens for efficient processing. Its architecture consists of a SigLIP-based vision encoder [64], a lightweight MLP projector, and a Qwen2-based language model [63]. For fine-tuning, NVILA applies lower learning rates to the vision encoder’s Layer-Norms while optimizing the language model. This allows robust AV adaptation under limited compute budgets [32].

LoRA (Low-Rank Adaptation) [22] offers an efficient alternative by focusing on specific parameters within the model. Only low-rank matrices are introduced in certain layers — primarily within the self-attention and feedforward blocks — allowing the majority of the pre-trained parameters to remain frozen. This technique reduces the memory and computational load while achieving task-specific adjustments by updating only the new, smaller matrices. By leveraging fine-tuning techniques for AV tasks and the efficiency gains of LoRA, exploring V^2LMs as a solution to AV perception challenges has the potential to address critical limitations of existing methods, enhancing robustness against adversarial attacks while avoiding the performance degradation common in traditional defenses.

Solo vs. Tandem Design Comparison. Deploying V^2LMs to enhance robustness against AEs in AVs raises important considerations for their implementation. Given the diverse range of tasks that AVs must perform, it is crucial to determine the best strategy for utilizing them. One approach involves using a separate V^2LM for each specific AV task to improve robustness within each module, ensuring specialized

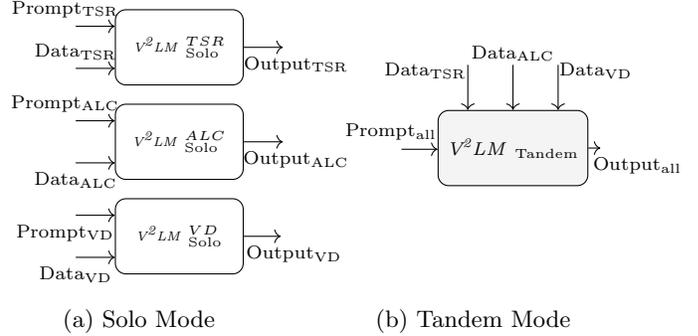


Figure 3: Comparison of Solo and Tandem Modes

and precise detection capabilities. This design, referred to as *Solo Mode*, is illustrated in Figure 3a, where individual V^2LMs are dedicated to tasks such as TSR, ALC, and VD. The formal representation of this design is:

$$Output_i = V^2LM_{solo}^i(Data_i, Prompt_i) \quad (4)$$

where $i \in \{TSR, ALC, VD\}$. Here, $Data_i$ denotes the dataset specific to each perception task T_i , and $Prompt_i$ is the prompt customized to fine-tune $V^2LM_{solo}^i$ for that particular task, ensuring task-specific optimization.

Alternatively, a single V^2LM can handle all AV tasks in a tandem approach, providing a unified method for improving robustness across multiple modules, as depicted in Figure 3b. In this design, multiple image-query pairs—each corresponding to a different task—are combined into a single input, where images are concatenated using a separator, and queries are merged in the same order. The model processes this structured batch simultaneously within one forward pass, extracting task-specific outputs for each image-query pair independently while treating the collection as a cohesive input during execution:

$$\{Output_{TSR}, Output_{ALC}, Output_{VD}\} = V^2LM_{tandem}(Data_{all}, Prompt_{all}) \quad (5)$$

where $Data_{all}$ represents the combined dataset of all tasks, and $Prompt_{all}$ includes the concatenated prompts corresponding to each task. Moreover, a single V^2LM could reduce memory and computational overhead, which is critical in the resource-constrained environment of AVs. In such systems, optimizing both efficiency and memory is vital, making a unified V^2LM a more practical solution.

V^2LM -Augmented Perception Systems. Integrating V^2LMs with AV perception systems offers the

potential to strengthen resilience against AEs in real-time. Currently, AV perception systems face significant limitations in mitigating attacks, which can lead to dangerous misinterpretations of sensor data. As Cao *et al.* [10] demonstrated, despite efforts to enhance the security of AV perception, significant vulnerabilities persist and traditional detection mechanisms often fail to mitigate these threats, leading to potentially dangerous consequences for AV decision-making and safety.

By integrating V^2LM in the end-to-end AV stack, AVs can benefit from its ability to process data, thereby working in parallel with perception tasks to enhance robustness and support accurate decision-making. The outputs of the perception system and V^2LM can then be used by the downstream *planning* and *control* modules to act accordingly in the presence of AEs. This context is depicted in Figure 4.

It is crucial to evaluate the latency of V^2LM integration to ensure it can enhance robustness against AEs in real-time without affecting the efficiency of the perception system. AV vision algorithms can process images at a rate of 2-4 frames per second (fps), which corresponds to approximately 250-500 ms per image [6]. As a result, V^2LM should process inputs within the required real-time constraints to support robust perception.

4 Experimental Evaluation

4.1 Experimental Setup

To evaluate the efficacy of V^2LMs , the focus was on three critical tasks: Traffic Sign Recognition (TSR), Automated Lane Centering (ALC) and Vehicle Detection (VD). For the former, the German Traffic Sign Recognition Benchmark (GTSRB) [54] was utilized, which includes images captured under various conditions such as different lighting and distances. Each picture from this dataset, which has 42 classes, was sent with the prompt *"Identify this traffic sign."* to the respective V^2LM . For the ALC task, a dataset was generated using the CARLA simulator [12], which includes 6,000 training and 2,000 testing images. These images, taken from the driver’s perspective under various weather conditions and times of day, are classified into three categories indicating the next move: *Straight*, *Left*, and *Right*. The prompt used for this task was, *"As a car driver, at which direction should you turn the steering wheel?"*

For VD, CARLA was also used to create 7,000 training and 3,000 testing images. Captured from different viewpoints, under diverse weather conditions, at various times of day, and from different distances,

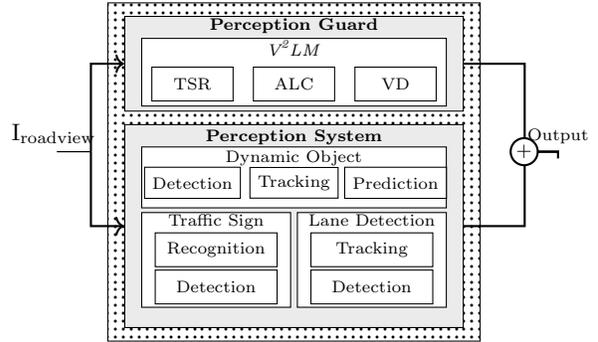


Figure 4: Possible integration of V^2LM with perception system for attack mitigation. $I_{roadview}$ is the image captured by AV cameras, with the output sent to the planning and control modules for action.

these images are categorized into *Car presence* and *Car absence*. If a car is present, the bounding box coordinates are provided with x and y representing the center, and H and W indicating the height and width of the box. The prompt used for this task was, *"If a car is detected, provide the center coordinates and the dimensions of the bounding box for the car"*.

The evaluation began with assessing the VLMs’ performance on zero-shot tasks across the test datasets. Subsequently, these VLMs were fine-tuned on AV-specific training data — yielding V^2LMs — with their performance re-evaluated on the same test datasets to ensure comparability. The objective is to use a V^2LM to enhance the ALC, TSR, and VD modules within the AV perception system against adversarial attacks. In the next step, AEs were generated from the same test datasets to assess robustness.

Although the fine-tuned models were trained on AV-specific data, the AEs remained *unseen* during training, allowing a reliable evaluation of each V^2LM ’s resilience to attacks on TSR, ALC, and VD. To achieve this, three types of black-box attacks were implemented. The first type of attack involves adversarial manipulation of traffic signs. In the study by Zhong *et al.* [69], shadows are utilized to conduct attacks on TSR algorithms, as shown in Figure 5a.

This method employs shadows as a non-invasive mechanism to create physical AEs. By optimizing shadow properties such as shape and opacity using a differentiable renderer, the technique manipulates images under black-box conditions to induce misclassifications. It achieves a success rate of 90.47% on the GTSRB dataset, demonstrating its effectiveness and highlighting the vulnerabilities of current detection systems to such subtle manipulations.

The second type of attack targets the ALC mechanism of AVs. In the study by Sato *et al.* [48], the Dirty Road Patch (DRP) attack framework specifi-

Table 2: Comparison of V^2LM Performance in AV Tasks. Green Color shows Improvement of Tandem over Solo Model. LLaVA-13B is abbreviated for LLaVA-13B-LoRA.

Task	Model	Accuracy			F1-Score			Precision			Recall		
		Solo	Tandem	Diff	Solo	Tandem	Diff	Solo	Tandem	Diff	Solo	Tandem	Diff
TSR	LLaVA-7B	95.93%	97.38%	1.45%	94.62%	91.22%	-3.40%	94.94%	91.08%	-3.86%	94.31%	91.37%	-2.94%
	LLaVA-13B	97.15%	98.14%	0.99%	96.09%	93.21%	-2.88%	96.25%	92.35%	-3.90%	95.94%	94.10%	-1.84%
	MoE-LLaVA	95.24%	96.91%	1.67%	94.73%	95.78%	1.05%	94.51%	95.26%	0.75%	94.96%	96.31%	1.35%
	MobileVLM	88.19%	90.34%	2.15%	87.68%	90.18%	2.50%	86.17%	90.73%	4.56%	89.25%	89.65%	0.40%
	Qwen-7B	95.80%	94.13%	-1.67%	95.78%	94.27%	-1.51%	95.96%	94.63%	-1.33%	95.64%	93.88%	-1.76%
ALC	NVILA-8B	99.04%	99.14%	0.1%	99.10%	99.20%	0.1%	99.11%	99.19%	0.08%	99.13%	99.14%	0.01%
	LLaVA-7B	98.11%	95.41%	-2.70%	98.05%	95.24%	-2.81%	98.20%	95.21%	-2.99%	97.91%	95.29%	-2.62%
	LLaVA-13B	99.30%	98.31%	-0.99%	99.46%	97.85%	-1.61%	99.46%	98.32%	-1.14%	99.47%	97.40%	-2.07%
	MoE-LLaVA	96.86%	92.53%	-4.33%	96.23%	91.83%	-4.40%	97.63%	92.37%	-5.26%	96.85%	91.30%	-5.55%
	MobileVLM	84.41%	86.61%	2.20%	87.68%	89.53%	1.85%	86.17%	88.93%	2.76%	89.26%	90.10%	0.84%
VD	Qwen-7B	93.75%	93.91%	0.16%	93.83%	93.76%	-0.07%	94.39%	94.48%	0.09%	93.27%	93.05%	-0.22%
	NVILA-8B	99.51%	99.41%	-0.1%	99.45%	99.52%	0.07%	99.34%	99.44%	0.1%	99.55%	99.61%	-0.11%
	LLaVA-7B	95.84%	91.90%	-3.94%	94.56%	92.82%	-1.74%	95.12%	94.71%	-0.41%	94.01%	91.21%	-2.80%
	LLaVA-13B	97.05%	98.96%	1.91%	96.76%	95.97%	-0.79%	97.33%	97.88%	0.55%	96.21%	94.14%	-2.07%
	MoE-LLaVA	95.52%	93.42%	-2.10%	94.78%	92.30%	-2.48%	95.31%	93.83%	-1.48%	94.27%	90.82%	-3.45%
VD	MobileVLM	86.61%	88.01%	1.40%	88.42%	89.38%	0.96%	86.76%	87.52%	0.76%	90.15%	91.34%	1.19%
	Qwen-7B	96.02%	95.01%	-1.01%	96.55%	96.17%	-0.38%	97.76%	97.35%	-0.41%	94.91%	95.01%	0.10%
	NVILA-8B	99.93%	99.95%	0.02%	95.90%	99.84%	-0.06%	99.81%	99.79%	-0.02%	99.88%	99.91%	0.1%



(a) Shadow Attack (b) DRP Attack (c) RAUCA
Figure 5: Examples of adversarial attacks targeting AV perception.

cally targets ALC systems in AVs, exploiting vulnerabilities in deep learning-based lane detection. This method employs an optimization-based approach to systematically generate these patches, as illustrated in Figure 5b, considering real-world conditions such as lighting and camera angles to ensure effectiveness across different environmental scenarios. The optimized DRPs cause the AV to make incorrect steering decisions, which were demonstrated to be highly successful in real-world driving scenarios with a success rate exceeding 97.5%

Figure 5c shows the third type of attack which focuses on the adversarial camouflage of vehicles. Zhou *et al.* [70] propose a physical adversarial attack known as the Robust and Accurate UV-map-based Camouflage Attack (RAUCA) to deceive VD algorithms such as YOLOv3 [45]. It employs a technique utilizing a differentiable neural renderer, which allows for the optimization of adversarial camouflages through gradient back-propagation, enhancing both the robustness and precision of the attacks under varying environmental conditions. Their method achieved an attack success rate of 97.48% on the target detection models, demonstrating the significant vulnerability of these systems to such sophisticated camouflage attacks.

4.2 RQ1: Fine-Tuning Increases Detection Performance

Table 1 shows VLMs’ zero-shot performance on the test dataset, which performed poorly in ALC and

TSR tasks, indicating difficulties in these specific AV applications. Although the models demonstrated decent performance in the VD task, this success may partly be due to their pre-training on large, diverse image datasets, which likely enhanced their general visual recognition capabilities. However, they struggled to detect objects in images taken during darkness, rain, or when the objects were far away, showing their difficulty in maintaining high performance under varying and challenging environmental conditions despite their general visual understanding capabilities.

To improve their performance, these models were fine-tuned using the same training datasets and prompts as before, and then their performance was re-evaluated with the same test dataset. Table 2 shows notable accuracy improvements for all tasks after fine-tuning. For ALC, accuracy increased from 20.71%–50.91% to 86.61%–99.51%. For TSR, accuracy saw a substantial rise from 1.19%–31.24% to 90.34%–99.14%. In the VD task, accuracy improved from 60.64%–92.06% to 88.01

These findings suggest that although they initially underperformed on AV tasks, fine-tuning them with relevant datasets can lead to substantial performance improvements, showing their potential utility in AV applications. Intersection over Union (IoU) measures the overlap between the predicted output and the ground truth in tasks like object detection. A higher IoU indicates more accurate localization, making it

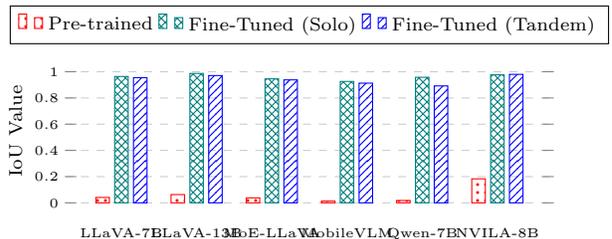


Figure 6: IoU values before and after fine-tuning.

a crucial metric for evaluating model performance in AV perception tasks. After fine-tuning, the IoU values improved significantly, with LLaVA-7B increasing from 0.041 to 0.964, LLaVA-13B-LoRA from 0.063 to 0.987, MoE-LLaVA from 0.038 to 0.946, MobileVLM from 0.013 to 0.926, Qwen-7B from 0.018 to 0.9578, and NVILA-8B from 0.183 to 0.9768, as shown in Figure 6.

4.3 RQ2: V^2LMs Demonstrate Robustness under Attacks

Adversarial attacks pose a serious threat to deep learning-based AV perception systems. Prior works have shown that DNN models suffer significant performance degradation across core AV tasks, including TSR, ALC, and VD, as summarized in Table 5. For example, the GTSRB-CNN model accuracy collapses to just 1.77% under the Shadow Attack [69], and even adversarial training—one of the most prominent defense methods—only modestly improves accuracy to 25.57%. Similarly, for ALC, OpenPilot-ALC performance degrades dramatically to 2.50% under the DRP Attack [48], with established defense strategies such as JPEG compression [20], Gaussian noise addition [61], and autoencoder-based denoising [19] achieving negligible improvements (around 3%). In VD, YOLOv3 experiences a drastic accuracy drop to 2.52% under RAUCA Attack [70], with no effective defense method proposed.

Building on these findings, we conduct our own robustness evaluation. To provide a fair and rigorous comparison, we evaluated traditional task-specific DNNs—YOLOv5-c1s [26, 44] (TSR), CLRerNet [21] (ALC), and YOLOv5-dt (VD)—on benign and previously *unseen* adversarial datasets (Table 3). Our results confirm severe performance degradation, with accuracy reductions of 39.06% for TSR, 40.40% for

Table 3: Performance of DNN models evaluated on benign (B) and unseen adversarial (A) datasets across AV tasks: TSR, ALC, and VD. The difference (Diff) indicates the accuracy degradation caused by adversarial attacks.

Task	Model	Type	Accuracy	F1-Score	Precision	Recall
TSR	Yo1ov5-c1s	A	55.82%	58.10%	60.56%	55.82%
		B	94.88%	93.44%	94.88%	93.92%
		Diff	-39.06%	-35.34%	-34.32%	-38.10%
ALC	CLRerNet	A	50.51%	44.86%	47.66%	50.51%
		B	90.91%	90.83%	91.44%	90.91%
		Diff	-40.40%	-45.97%	-43.78%	-40.40%
VD	Yo1ov5-dt	A	62.27%	57.44%	51.37%	64.11%
		B	97.01%	96.60%	96.66%	97.01%
		Diff	-34.74%	-39.16%	-45.29%	-32.91%

ALC, and 32.91% for VD, clearly highlighting their vulnerability to novel attacks.

In contrast, our evaluations of V^2LMs under the same unseen adversarial conditions reveal inherently superior robustness. Specifically, V^2LMs experienced substantially smaller accuracy drops: only 8.62%–15.81% for TSR, 4.94%–14.32% for ALC, and 5.23%–7.06% for VD, consistently maintaining high adversarial accuracy without additional defenses. These results strongly support the potential of V^2LMs to enhance robustness and reliability of AV perception systems.

4.4 RQ3: Tandem V^2LMs Provide Similar Performance at Lower Memory Footprint

This section evaluates whether a single *tandem* V^2LM can handle multiple defense tasks in detecting AEs across AV systems, compared to using separate *solo* V^2LMs for each task. To assess this, VLMs are fine-tuned for three tasks — TSR, ALC, and VD — simultaneously using the same prompts. After fine-tuning, both solo and tandem V^2LMs were then evaluated on each task to determine its performance on the same test data.

Table 2 presents the evaluation results for the tandem design in non-adversarial scenarios, where a single V^2LM was trained to handle all three tasks simultaneously. The results reveal that the tandem design achieves high accuracy across perception tasks, often matching or surpassing the performance of solo models. This shows the effectiveness of the tandem design in maintaining robust performance across diverse tasks. Table 4 further demonstrates the resilience of V^2LMs in AV tasks under adversarial conditions.

MobileVLM, MoE-LLaVA, LLaVA-7B, NVILA-8B, Qwen-VL, and LLaVA-13B-LoRA allocated 6.03GB, 11.21GB, 13.56GB, 15.23GB, 16.58GB, and 26.15GB of storage, respectively. These results suggest that a single tandem V^2LM can generalize well across multiple tasks, providing robust performance comparable to the solo design, which requires 3x more storage for separate models. The tandem design offers a significant advantage in efficiency by maintaining similar performance while requiring much less storage.

5 Discussion

AV perception systems typically target latencies below 100 milliseconds to meet real-time operational requirements, especially in high-speed driving contexts [25]. In our experiments, the inference time

Table 4: Comparison of V^2LM Performance against AEs. Green Color shows Improvement of Tandem over Solo Model. LLaVA-13B is abbreviated for LLaVA-13B-LoRA.

Task	Model	Accuracy			F1-Score			Precision			Recall		
		Solo	Tandem	Diff	Solo	Tandem	Diff	Solo	Tandem	Diff	Solo	Tandem	Diff
TSR	LLaVA-7B	80.12%	86.51%	6.39%	86.44%	85.89%	-0.55%	85.24%	86.12%	0.88%	87.69%	85.67%	2.02%
	LLaVA-13B	86.53%	89.01%	2.48%	86.08%	88.95%	2.87%	85.13%	89.31%	4.18%	87.06%	88.60%	1.54%
	MoE-LLaVA	80.03%	81.52%	1.49%	80.75%	82.34%	1.59%	82.13%	83.26%	1.13%	79.43%	81.45%	2.02%
	MobileVLM	79.57%	79.13%	-0.44%	77.55%	79.88%	2.33%	78.02%	80.75%	2.73%	77.10%	79.04%	1.94%
	Qwen-7B	82.83%	81.66%	-1.17%	83.16%	83.43%	0.27%	85.73%	86.52%	0.79%	80.57%	80.77%	0.20%
	NVILA-8B	86.04%	X%	-X%	87.15%	X%	-X%	90.28%	X%	-X%	86.09%	X%	-X%
ALC	LLaVA-7B	86.07%	84.83%	-1.24%	87.01%	85.38%	-1.63%	86.68%	84.65%	-2.03%	87.53%	86.13%	-1.40%
	LLaVA-13B	90.05%	86.69%	-3.36%	90.69%	87.38%	-3.31%	89.67%	86.34%	-3.33%	91.75%	88.46%	-3.29%
	MoE-LLaVA	82.54%	82.26%	-0.28%	83.27%	83.06%	-0.21%	84.46%	82.49%	-1.97%	82.13%	83.65%	1.52%
	MobileVLM	78.21%	78.23%	0.02%	80.92%	81.21%	0.29%	83.16%	81.96%	-1.20%	78.80%	80.49%	1.69%
	Qwen-7B	88.81%	86.57%	-2.24%	89.48%	85.48%	-4%	91.06%	89.61%	-1.45%	88.01%	86.57%	-1.44%
	NVILA-8B	99.25%	X%	-X%	99.13%	X%	-X%	99.01%	X%	-X%	99.00%	X%	-X%
VD	LLaVA-7B	89.23%	88.02%	-1.21%	87.84%	86.59%	-1.25%	89.06%	87.89%	-1.17%	86.67%	85.33%	-1.34%
	LLaVA-13B	91.82%	90.09%	-1.73%	90.72%	90.41%	-0.31%	90.10%	91.03%	0.93%	91.35%	89.80%	-1.55%
	MoE-LLaVA	88.46%	86.54%	-1.92%	87.16%	86.21%	-0.95%	88.28%	86.77%	-1.51%	86.07%	85.67%	-0.40%
	MobileVLM	80.13%	80.15%	0.02%	82.88%	82.97%	0.09%	84.14%	83.81%	-0.33%	81.67%	82.16%	0.49%
	Qwen-7B	89.05%	89.06%	0.01%	90.15%	90.57%	0.42%	92.83%	92.77%	-0.06%	87.60%	88.46%	0.86%
	NVILA-8B	99.81%	X%	-X%	99.82%	X%	-X%	99.87%	X%	-X%	99.81%	X%	-X%

t_{V^2LM} for LLaVA-7B—a 2023 model—was measured at 851 ms on an NVIDIA A100 GPU with 40 GB of VRAM [1], clearly exceeding the acceptable threshold t_{PS} for AV deployment. In contrast, the release of NVILA in late 2024 marked a significant improvement, reducing t_{V^2LM} to just 80 ms. This 10× reduction highlights the rapid evolution of VLMs toward real-time readiness in AV perception pipelines.

Despite this progress, deploying VLMs on embedded AV hardware remains challenging due to limitations in compute power and energy efficiency. High-performance models such as LLaVA-7B and NVILA, though effective on server-grade GPUs, often require substantial memory and parallel processing capabilities that are impractical for in-vehicle deployment. One potential solution to reduce hardware demands is quantization—a widely used model compression technique for LLMs that improves computational efficiency by converting high-precision data types to lower-precision formats [23]. This process significantly reduces memory usage and model size, making it more feasible to run VLMs on edge devices. However, it may also introduce quantization errors, which could degrade precision [29]. We applied quantization to LLaVA-7B, expecting lower resource usage. Although initially expected to reduce latency, the quantization approach did not yield the desired outcome; instead, when tested on the NVIDIA A100 40 GB [1], it exhibited latencies ranging from 0.851 s

to 2.158 s and 11.657 s at full precision, 8-bit, and 4-bit levels, respectively. Based on our preliminary analysis, we found that the main source of this issue could be due to an increase in demand stemming from the additional operations required to maintain accuracy in image processing tasks.

Another complementary strategy is adaptive image downsampling based on scene complexity. For example, in visually simple highway scenes, input images could be processed at lower resolutions to reduce computational load, while complex urban scenes could use higher resolutions to preserve important details. This approach creates a trade-off between efficiency and detection quality, as previously observed in object detection for AVs [7]. By adjusting the input resolution based on the type of scene, AV systems can improve resource management while maintaining reliable performance across different driving environments. In parallel, compression techniques such as model pruning and knowledge distillation (KD) have been successfully applied to reduce model size and inference cost in vision tasks. For instance, Rosa *et al.* [47] demonstrate that KD can yield a student model that is over 90% smaller than its teacher network while maintaining competitive accuracy in semantic segmentation. These techniques can complement downsampling strategies, further improving the feasibility of deploying VLMs on resource-constrained AV hardware.

Table 5: Traditional Models Accuracy under Adversarial Attacks and Defense Strategies.

Task	Model	Attack Type	Under-Attack Accuracy	Defense Strategy	Post-Defense Accuracy
TSR	GTSRB-CNN	Shadow [69]	1.77%	Adversarial Training [36]	25.57%
ALC	OpenPilot-ALC	DRP [48]	2.50%	JPEG Compression [20]	≈3% (no effective improvement)
				Bit-Depth Reduction [62]	
				Gaussian Noise [61] Median Blurring [62] Autoencoder [19]	
VD	YOLOv3	RAUCA [70]	2.52%	No Defense Proposed	N/A

6 Conclusion

In this work, we present Vehicle Vision Language Models (V^2LMs) as a novel approach to enhance the robustness of AV perception systems against adversarial attacks. V^2LMs outperform task-specific DNNs by maintaining high adversarial accuracy without requiring adversarial training, addressing a key limitation of traditional defense methods. The study

evaluates two deployment strategies: *Solo Mode*, where separate VLMs are fine-tuned for common AV perception tasks such as TSR, ALC, and VD, as well as *Tandem Mode*, which uses a single unified V^2LM for all three tasks, reducing storage requirements while maintaining comparable performance. Experimental results demonstrate that traditional DNN models suffer performance drops of 33%–46% under adversarial attacks, whereas V^2LMs maintain significantly higher accuracy, with reductions of less than 8% on average. Additionally, this work explores the feasibility of integrating V^2LMs into AV systems, proposing solutions such as model compression and adaptive image downsampling to mitigate latency challenges. This paper shows that V^2LMs significantly enhance AV perception robustness, offering a promising path toward safer autonomous driving.

References

- [1] NVIDIA A100 Tensor Core GPU. <https://www.nvidia.com/en-us/data-center/a100/>. Accessed: 2023-11-14. 9
- [2] Mohammed Aldeen, Pedram MohajerAnsari, Jin Ma, Mashrur Chowdhury, Long Cheng, and Mert D. Pesé. An initial exploration of employing large multimodal models in defending against autonomous vehicles attacks. In *2024 IEEE Intelligent Vehicles Symposium (IV)*, pages 3334–3341, 2024. 3
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 2, 5
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 1(8), 2023. 2
- [5] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1728–1738, 2021. 4
- [6] Krzysztof Blachut, Michal Danilowicz, Hubert Szolc, Mateusz Wasala, Tomasz Kryjak, and Mateusz Komorkiewicz. Automotive perception system evaluation with reference data from a uav’s camera using aruco markers and dcnn. *Journal of Signal Processing Systems*, 94(7):675–692, 2022. 6
- [7] Imene Bouderbail, Abdenour Amamra, and Mohamed Akrem Benatia. How would image downsampling and compression impact object detection in the context of self-driving vehicles? In *International Conference on Computing Systems and Applications*, pages 25–37. Springer, 2020. 9
- [8] Michael D Breitenstein, Fabian Reichlin, Bastian Leibe, Esther Koller-Meier, and Luc Van Gool. Online multiperson tracking-by-detection from a single, uncalibrated camera. *IEEE transactions on pattern analysis and machine intelligence*, 33(9):1820–1833, 2010. 4
- [9] Yulong Cao, Ningfei Wang, Chaowei Xiao, Dawei Yang, Jin Fang, Ruigang Yang, Qi Alfred Chen, Mingyan Liu, and Bo Li. Invisible for both camera and lidar: Security of multi-sensor fusion based perception in autonomous driving under physical-world attacks. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 176–194. IEEE, 2021. 3
- [10] Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Won Park, Sara Rampazzi, Qi Alfred Chen, Kevin Fu, and Z Morley Mao. Adversarial sensor attack on lidar-based perception in autonomous driving. In *Proceedings of the 2019 ACM SIGSAC conference on computer and communications security*, pages 2267–2281, 2019. 6
- [11] Claudio Caraffi, Tomáš Vojtř, Jiří Trefný, Jan Šochman, and Jiří Matas. A system for real-time detection and tracking of vehicles from a single car-mounted camera. In *2012 15th international IEEE conference on intelligent transportation systems*, pages 975–982. IEEE, 2012. 4
- [12] CARLA. Carla simulator, 2023. 6
- [13] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2818–2829, 2023. 5
- [14] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. 5
- [15] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023. 2
- [16] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran Wang. Receive, reason, and react: Drive as you say, with large language models in autonomous vehicles. *IEEE Intelligent Transportation Systems Magazine*, 2024. 3
- [17] Yogita Dubey, Yashraj Tarte, Nikhil Talatule, Bhargav Sable, Chetan Taywade, and Roshan Umate. An artificial intelligence based autonomous road lane detection and navigation system for vehicles. 2024. 4

- [18] Kevin Eykholt, Ivan Evtimov, Earleence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1625–1634, 2018. 3
- [19] Shixiang Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial examples. *arXiv preprint arXiv:1412.5068*, 2014. 8, 9
- [20] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens Van Der Maaten. Countering adversarial images using input transformations. *arXiv preprint arXiv:1711.00117*, 2017. 2, 4, 8, 9
- [21] Hiroto Honda and Yusuke Uchida. CLRerNet: Improving Confidence of Lane Detection with LaneIoU. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2024. 2, 8
- [22] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 5
- [23] IBM. What is quantization?, 2023. Accessed: 2023-11-14. 9
- [24] Wei Jia, Zhaojun Lu, Haichun Zhang, Zhenglin Liu, Jie Wang, and Gang Qu. Fooling the eyes of autonomous vehicles: Robust physical adversarial examples against traffic sign recognition systems. *arXiv preprint arXiv:2201.06192*, 2022. 1
- [25] Tianshi Jin, Weiping Ding, Mingliang Yang, Honglin Zhu, and Peisong Dai. Benchmarking perception to streaming inputs in vision-centric autonomous driving. *Mathematics*, 11(24):4976, 2023. 8
- [26] Glenn Jocher, Ayush Chaurasia, Alex Stoken, Jirka Borovec, et al. YOLOv5: PyTorch implementation of YOLO object detector. <https://github.com/ultralytics/yolov5>, 2020. 2, 8
- [27] Zelun Kong, Junfeng Guo, Ang Li, and Cong Liu. Physgan: Generating physical-world-resilient adversarial examples for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14254–14263, 2020. 3
- [28] Bin Lin, Zhenyu Tang, Yang Ye, Jiayi Cui, Bin Zhu, Peng Jin, Junwu Zhang, Munan Ning, and Li Yuan. Moe-llava: Mixture of experts for large vision-language models. *arXiv preprint arXiv:2401.15947*, 2024. 2
- [29] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024. 9
- [30] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 2
- [31] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 5
- [32] Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, et al. Nvila: Efficient frontier visual language models. *arXiv preprint arXiv:2412.04468*, 2024. 2, 5
- [33] Zhou Lu. A theory of multimodal learning. *Advances in Neural Information Processing Systems*, 36, 2024. 4
- [34] Chen Ma, Ningfei Wang, Qi Alfred Chen, and Chao Shen. Slowtrack: Increasing the latency of camera-based perception in autonomous driving using adversarial examples. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4062–4070, 2024. 2
- [35] Chen Ma, Ningfei Wang, Zhengyu Zhao, Qian Wang, Qi Alfred Chen, and Chao Shen. Controlloc: Physical-world hijacking attack on visual perception in autonomous driving. *arXiv preprint arXiv:2406.05810*, 2024. 2
- [36] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017. 4, 9
- [37] Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. Gpt-driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023. 3
- [38] Pedram MohajerAnsari, Alkim Domeke, Jan de Voor, Arkajyoti Mitra, Grace Johnson, Amir Salarpour, Habeeb Olufowobi, Mohammad Hamad, and Mert D Pesé. Discovering new shadow patterns for black-box attacks on lane detection of autonomous vehicles. *arXiv preprint arXiv:2409.18248*, 2024. 1
- [39] Pedram MohajerAnsari, Amir Salarpour, David Fernandez, Cigdem Kokenoz, Bing Li, and Mert D Pesé. Attention-aware temporal adversarial shadows on traffic sign sequences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025. 1
- [40] Mohammad Omama, Pranav Inani, Pranjal Paul, Sarat Chandra Yellapragada, Krishna Murthy Jatavallabhula, Sandeep Chinchali, and Madhava Krishna. Alt-pilot: Autonomous navigation with language augmented topometric maps. *arXiv preprint arXiv:2310.02324*, 2023. 3
- [41] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense

- to adversarial perturbations against deep neural networks. In *2016 IEEE symposium on security and privacy (SP)*, pages 582–597. IEEE, 2016. 2, 4
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5
- [43] Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John C Duchi, and Percy Liang. Adversarial training can hurt generalization. *arXiv preprint arXiv:1906.06032*, 2019. 2
- [44] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016. 2, 8
- [45] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 7
- [46] Md Kaykobad Reza, Ashley Prater-Bennette, and M Salman Asif. Robust multimodal learning with missing modalities via parameter-efficient adaptation. *arXiv preprint arXiv:2310.03986*, 2023. 4
- [47] Ciro Rosa and Nina Hirata. Knowledge distillation for reduced footprint semantic segmentation with the u-net. In *Proceedings of the 40th ACM/SIGAPP Symposium on Applied Computing*, pages 655–662, 2025. 9
- [48] Takami Sato, Junjie Shen, Ningfei Wang, Yunhan Jia, Xue Lin, and Qi Alfred Chen. Dirty road can attack: Security of deep learning based automated lane centering under {Physical-World} attack. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 3309–3326, 2021. 2, 6, 8, 9
- [49] Wilko Schwarting, Javier Alonso-Mora, and Daniela Rus. Planning and decision-making for autonomous vehicles. *Annual Review of Control, Robotics, and Autonomous Systems*, 1:187–210, 2018. 1
- [50] Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, and Mingyu Ding. Languagempc: Large language models as decision makers for autonomous driving. *arXiv preprint arXiv:2310.03026*, 2023. 3
- [51] Junjie Shen, Ningfei Wang, Ziwen Wan, Yunpeng Luo, Takami Sato, Zhisheng Hu, Xinyang Zhang, Shengjian Guo, Zhenyu Zhong, Kang Li, et al. Sok: On the semantic ai security in autonomous driving. *arXiv preprint arXiv:2203.05314*, 2022. 3
- [52] Ruoyu Song, Muslum Ozgur Ozmen, Hyungsub Kim, Raymond Muller, Z Berkay Celik, and Antonio Bianchi. Discovering adversarial driving maneuvers against autonomous vehicles. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 2957–2974, 2023. 1
- [53] NN Sriram, Tirth Maniar, Jayaganesh Kalyanasundaram, Vineet Gandhi, Brojeshwar Bhowmick, and K Madhava Krishna. Talk to the vehicle: Language conditioned autonomous navigation of self driving cars. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5284–5290. IEEE, 2019. 3
- [54] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, (0):–, 2012. 6
- [55] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 3
- [56] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2
- [57] Nesrine Triki, Mohamed Karray, and Mohamed Ksantini. A real-time traffic sign recognition method using a new attention-based deep convolutional neural network for smart vehicles. *Applied Sciences*, 13(8):4793, 2023. 4
- [58] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018. 2
- [59] Li-Hua Wen and Kang-Hyun Jo. Deep learning-based perception systems for autonomous driving: A comprehensive survey. *Neurocomputing*, 489:255–270, 2022. 1
- [60] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *International conference on machine learning*, pages 5286–5295. PMLR, 2018. 2, 4
- [61] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. *arXiv preprint arXiv:1711.01991*, 2017. 8, 9
- [62] Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. *arXiv preprint arXiv:1704.01155*, 2017. 9
- [63] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024. 5

- [64] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. [5](#)
- [65] Ce Zhang, Azim Eskandarian, and Xuelai Du. Attention-based neural network for driving environment complexity perception. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pages 2781–2787. IEEE, 2021. [4](#)
- [66] Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. [4](#)
- [67] Qingzhao Zhang, Shengtuo Hu, Jiachen Sun, Qi Alfred Chen, and Z Morley Mao. On adversarial robustness of trajectory prediction for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15159–15168, 2022. [1](#)
- [68] Zheng Zhao, Emilio Monti, Jens Lehmann, and Haytham Assem. Enhancing contextual understanding in large language models through contrastive decoding. *arXiv preprint arXiv:2405.02750*, 2024. [4](#)
- [69] Yiqi Zhong, Xianming Liu, Deming Zhai, Junjun Jiang, and Xiangyang Ji. Shadows can be dangerous: Stealthy and effective physical-world adversarial attack by natural phenomenon. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15345–15354, 2022. [1](#), [2](#), [6](#), [8](#), [9](#)
- [70] Jiawei Zhou, Linye Lyu, Daojing He, and Yu Li. Rauca: A novel physical adversarial attack on vehicle detectors via robust and accurate camouflage generation. *arXiv preprint arXiv:2402.15853*, 2024. [2](#), [7](#), [8](#), [9](#)