

Memory Makes The Poisons: Understanding and Mitigating Data Poisoning in LVLMs

Anonymous authors

Paper under double-blind review

Abstract

Large Vision-Language Models (LVLMs) are increasingly deployed in high-stakes applications, yet their training-time security remains poorly understood. As a prominent data poisoning attack specifically designed for LVLMs, ShadowCast Xu et al. (2024) achieves significant success in inducing targeted hallucinations, posing a serious threat to LVLM safety. ShadowCast’s success has been attributed to injected visual perturbations. Consequently, subsequent defenses have focused on visual purification; however, their effectiveness remains limited. **In this paper**, we present a re-analysis of the ShadowCast mechanism. Our key finding is that memorization during LVLM fine-tuning is an overlooked but major contributor to attack success, and it dominates at higher poison ratios. This factor has been largely overlooked in previous work. We further show that multimodal training exacerbates this vulnerability compared to unimodal settings. This insight fundamentally reframes both the threat model and the defense objective: if memorization is a major contributor, purification-only defenses are inherently insufficient in multimodal regimes. Motivated by this perspective, we propose *RejectShield*, a rejection-based defense that filters suspicious training samples prior to fine-tuning. Across extensive evaluations spanning 4 attack goals, 3 LVLMs, black-box and white-box attack settings, and 3 poisonings, RejectShield reduces the attack success rate by up to 99% while largely preserving model utility, significantly advancing defense effectiveness against LVLM poisoning. **Code, checkpoints and additional results are provided in the Supp.**

1 Introduction

Large Vision-Language Models (LVLMs) have recently achieved strong performance across a wide range of applications such as Visual Question Answering (Li et al., 2023b; Wang et al., 2025), Image Captioning (Alayrac et al., 2022; Zhang et al., 2026), Embodied AI / Robotics (Shen et al., 2023; Ju et al., 2026), Medical Imaging Assistants (Mullick et al., 2023; Xie et al., 2025; Liu et al., 2026), and Interactive Agents and Game Environments (Xu et al., 2023; Peng et al., 2026). Their ability to align visual and textual inputs makes them attractive for real-world applications, but it also raises safety and security concerns (Zhao et al., 2023; Nguyen et al., 2026; Ding et al., 2025; Lee et al., 2025). To improve downstream performance, LVLMs are often fine-tuned on crowd-sourced datasets. During fine-tuning, these models can be vulnerable to malicious data poisoning attacks that induce targeted hallucinations at deployment (Xu et al., 2024). Ensuring robust and responsible deployment of LVLMs therefore requires a deeper understanding of these vulnerabilities and defenses that go beyond straightforward adaptations of unimodal techniques.

Research Gap. Despite growing interest in LVLM safety and security, their vulnerability to data poisoning remains inadequately explored. A prominent data poisoning attack specifically designed to target LVLMs is ShadowCast (Xu et al., 2024), whose main idea is to introduce carefully-crafted adversarial visual perturbations into fine-tuning data to induce targeted hallucinations. ShadowCast has achieved significant attack success, posing a serious threat to LVLM safety. The effectiveness of ShadowCast has been primarily attributed to the injection of adversarial visual perturbations (Xu et al., 2024). Consequently, subsequent defenses have almost exclusively focused on purification approaches that aim to remove these perturbations. However, their defense effectiveness remains limited. This raises a fundamental question: *Is adversarial visual perturbation*

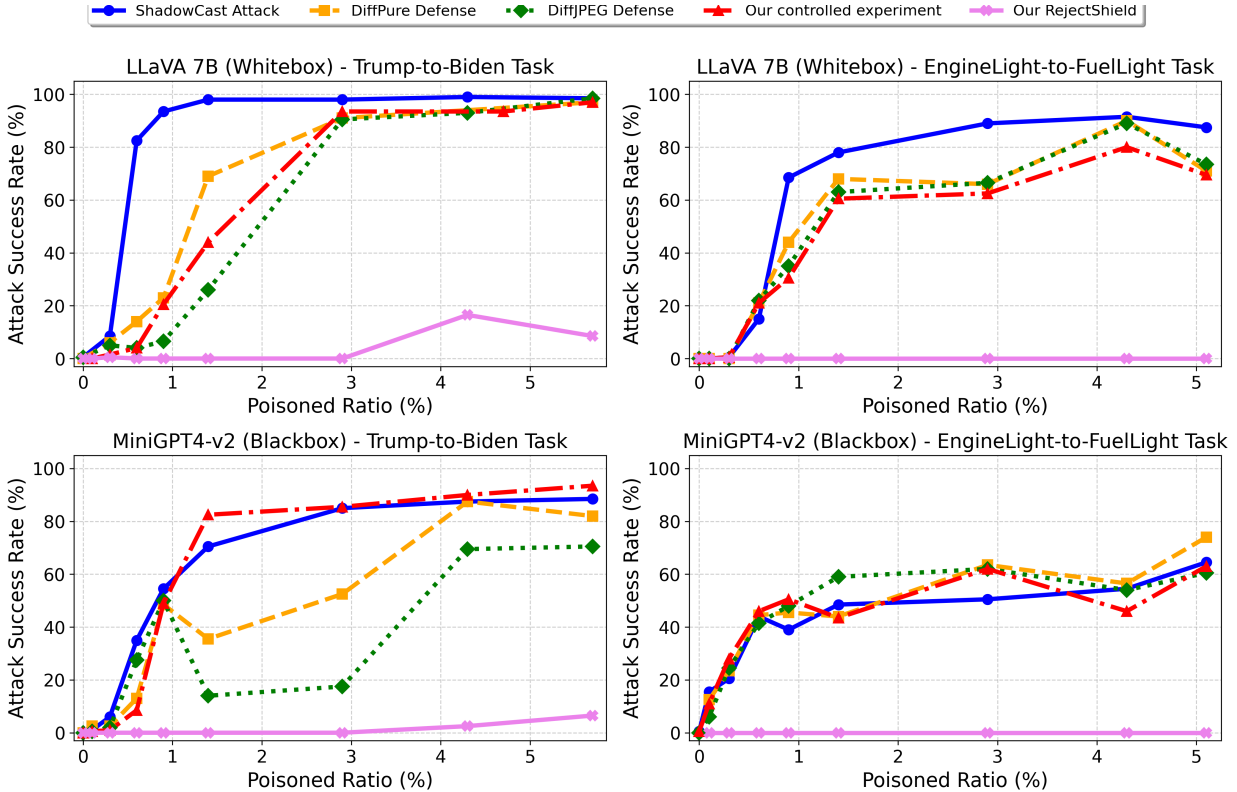


Figure 1: **Overview of our analysis of LVLm data poisoning via ShadowCast (Xu et al., 2024).** In the main paper, we analyze the standard ShadowCast setup (blue) using two primary LVLms (LLaVA 1.5 and MiniGPT4-v2 (Zhu et al., 2023; Liu et al., 2023)) on two Label Attack tasks (Trump-to-Biden and Engine-to-Fuel), under both white-box and black-box settings. (LLaVA-NeXT result is also included in Fig. 5 to test transfer to a newer LVLm). We report three observations. **(1) Existing purification-based defenses have limited effectiveness in this setting.** DiffPure (orange) (Nie et al., 2022) and DiffJPEG (green) (Shin et al., 2017) can help at low poison ratios, but their effectiveness drops as poison ratio increases. **(2) Data memorization is a major contributor to ShadowCast effectiveness in our settings and can dominate at higher poison ratios.** In controlled experiments (red) where we remove visual perturbations but keep the setup otherwise unchanged, attack success remains high, especially at higher poison ratios. Visual perturbations still provide efficiency gains at low poison ratios. This helps explain why defenses that only denoise visual perturbations are insufficient. **(3) RejectShield improves robustness.** Based on this analysis, we introduce RejectShield (pink), a rejection-based defense that reduces attack success rates by up to 99% while largely preserving model utility in our evaluated settings. Additional experiments are provided in the Supp.

the sole major factor underlying the vulnerability of LVLms to ShadowCast, or is there a deeper, previously overlooked mechanism at play?

In this paper, we present an in-depth study of LVLm poisoning attacks through ShadowCast. To examine whether adversarial visual perturbations are the sole major factor or if other mechanisms are involved, we design controlled experiments that isolate the effect of these perturbations. Through these experiments, we uncover that *data memorization is a previously overlooked but major contributor to attack success*, and it dominates at higher poison ratios. This factor has not been identified in prior work in data poisoning targeting LVLms. We further show that multimodal training in LVLms exacerbates this vulnerability compared to a comparable unimodal setting. This insight reframes both the threat model and the defense objective: defenses based solely on purification are insufficient in regimes where memorization plays a significant role. Based on

this finding, we introduce a rejection-based approach, *RejectShield*, to improve defense effectiveness against data poisoning attacks on LVLMs.

Our main contributions are:

- **Data memorization as a major contributor to LVLM poisoning attacks.** We design controlled experiments to examine whether adversarial visual perturbations are the sole major factor for the success of LVLM data poisoning, or if other mechanisms are involved. Our experiments provide evidence that data memorization during LVLM fine-tuning is an overlooked, major contributor to attack success, and it dominates at higher poison ratios. We further show in carefully-designed LVLM-vs-LLM experiments that multimodal inputs exacerbate this memorization vulnerability. This analysis explains why previous defenses based on purification often struggle against LVLM poisoning. See Sec. 3
- **A rejection-based defense, *RejectShield*.** Our new insight fundamentally reframes both the threat model and the defense objective: if memorization is a major contributor, purification-only defenses are inherently insufficient in multimodal regimes. Motivated by this perspective, we propose *RejectShield*, a simple rejection-based defense that filters suspected poisoned samples before LVLM fine-tuning. See Sec. 4
- **Extensive evaluation.** We evaluate *RejectShield* across extensive evaluations spanning 4 attack goals, 3 LVLMs, black-box and white-box attack settings, and 3 poisonings. In our evaluated settings, *RejectShield* reduces attack success rate by up to 99% while largely maintaining model utility. This achieves state-of-the-art defense effectiveness against LVLM poisoning.

2 Related Work

Large Vision Language Models (LVLMs). Large Vision-Language Models (LVLMs) extend Large Language Models (LLMs) by integrating vision encoders to handle both visual and textual modalities, excelling in tasks such as visual question answering, image captioning, multimodal dialogue, and robotics. They typically consist of a visual encoder (e.g., CLIP (Radford et al., 2021), ViT (Dosovitskiy et al., 2021)) and a large language model (e.g., LLaMA (Touvron et al., 2023), Vicuna (Chiang et al., 2023)), with modality fusion achieved through projection layers (Liu et al., 2023; Zhu et al., 2023; Li et al., 2023b) or attention mechanisms (Li et al., 2022; Alayrac et al., 2022; Li et al., 2023a). BLIP (Li et al., 2022) and BLIP-2 (Li et al., 2023b) introduce two-stage pipelines combining vision-language contrastive learning and instruction tuning. LLaVA (Liu et al., 2023) efficiently aligns CLIP features with LLMs using a lightweight projection head and has become one of the most widely adopted open-source LVLMs. MiniGPT-4 (Zhu et al., 2023), Otter (Li et al., 2023a), and InternGPT (Li et al., 2023a) further improve LVLM capabilities with improved visual grounding, instruction-following, and multilingual support, respectively. The development of LVLMs typically on a downstream task follows a two-stage process: (1) pre-training on large-scale datasets for general multimodal understanding, and (2) fine-tuning for task-specific alignment or instruction following. *In this work, we focus on the analysis of vulnerabilities of LVLMs during the fine-tuning stage.*

Data Poisoning Attack in LVLMs. Data poisoning attacks aim to inject malicious data during training to compromise a model’s behavior at inference (Steinhardt et al., 2017; Gu et al., 2017; Zhu et al., 2019). Prior research has primarily focused on unimodal settings, such as vision-only (Shafahi et al., 2018; Zhao et al., 2020; Turner et al., 2019) or text-only models (Wallace et al., 2021; Kurita et al., 2020). Recent benchmarking efforts such as PoisonBench (Fu et al., 2024) further emphasize the importance of systematic evaluation protocols for poisoning vulnerability in large language models. In multimodal settings, recent work explores additional poisoning vectors in modern LVLM pipelines, including knowledge poisoning and stealthy poisoning attacks in retrieval-augmented LVLMs (Zhang et al., 2025; Yu et al., 2025), as well as semantic-manipulation backdoor attacks (Zhong et al., 2025). These studies highlight that poisoning risks can arise through multiple mechanisms, including perturbation-based, retrieval-knowledge, and semantic backdoor pathways. Among these attacks, ShadowCast (Xu et al., 2024) is especially important because it targets the fine-tuning stage,

which is the standard step for adapting LVLMs to downstream tasks. ShadowCast injects imperceptible visual perturbations into fine-tuning data and induces targeted hallucinations in the resulting model, making it a practical and high-impact threat. Despite its effectiveness, the mechanism behind ShadowCast remains insufficiently understood, and prior explanations mainly attribute success to visual adversarial perturbations. As a result, follow-up defenses primarily focus on purifying visual perturbations, yet they often have limited effectiveness in this setting (Xu et al., 2024). In this work, we therefore focus on ShadowCast and provide evidence that data memorization during LVLm fine-tuning is a major contributor to attack success in our settings, especially at higher poison ratios, while visual perturbations remain important at low poison ratios. Building on this insight, we propose a rejection-based defense that substantially mitigates the attack while largely preserving model utility.

Memorization in Deep Neural Networks. While deep neural networks (DNNs) are generally expected to learn patterns from training data, prior studies show that they can also memorize training examples, even when labels are random (Zhang et al., 2017). This suggests that DNNs may encode sample-specific artifacts rather than only generalizable structure. In practice, this behavior can create security and reliability risks. For Large Language Models (LLMs), previous works show that models can memorize and leak rare or sensitive information from their training corpora (Carlini et al., 2021; Zhang et al., 2021; Lee et al., 2022). Extending these concerns to multimodal learning, (Jayaraman et al., 2024) explores memorization in contrastive vision-language models such as CLIP, showing recall of fine-grained visual details that are absent from paired captions. However, that analysis focuses on retrieval rather than generative LVLm behavior. In this work, we investigate data memorization in generative LVLms and identify a related failure mode: hallucinated outputs that appear to rely on memorized fine-tuning patterns rather than grounding in the current multimodal input (e.g., image and question). This motivates further work on improving grounding and factual consistency during LVLm fine-tuning.

3 Data Memorization during LVLm Fine-tuning

LVLms are vulnerable to data poisoning attacks Xu et al. (2024), but the mechanism behind attack success is still not fully understood. In Sec. 3.1, we review background and motivation using ShadowCast Xu et al. (2024). In Sec. 3.2, we investigate the mechanism behind ShadowCast and provide evidence that data memorization is a major contributor to attack success in our settings, especially at higher poison ratios, while visual perturbations remain important for attack efficiency at low poison ratios. In Sec. 3.3, we further compare unimodal and multimodal settings and show that multimodal inputs can exacerbate memorization effects in LVLms. Overall, these results provide a more complete picture of LVLm data poisoning behavior. Additional results and detailed experimental designs are provided in the Supp.

3.1 Background and Motivation

Background. Data poisoning attacks aim to inject malicious training samples into a model’s dataset to induce incorrect or attacker-controlled behavior at inference time Steinhardt et al. (2017); Gu et al. (2017); Zhu et al. (2019). In unimodal settings, extensive research has targeted vision models using clean-label Shafahi et al. (2018); Turner et al. (2019) and optimization-based poisons Geiping et al. (2021), as well as NLP models via weight-poisoning or trigger-based techniques Kurita et al. (2020); Wallace et al. (2021). However, in multimodal settings such as LVLms, data poisoning attacks are underexplored.

A recent pioneering study, ShadowCast Xu et al. (2024), exposes a novel threat to LVLms during the fine-tuning phase, a critical stage for adapting pre-trained LVLms to downstream tasks. Particularly, LVLms are fine-tuned with clean data $\mathcal{D}_{\text{clean}}$ including clean text-image pairs (x_c, y_c) .

ShadowCast induces targeted hallucinations by injecting carefully crafted poisoned training pairs (x_p, y_d) , referred to as *visually matching poison samples*. Each poisoned image x_p is optimized to be visually similar to x_d (representing the destination concept C_d) to humans, while also being similar in the LVLms visual latent feature space to an image x_o (representing the source concept C_o). The poison image x_p is then paired with caption y_d , which is the caption of x_d , forming a poisoned training pair (x_p, y_d) . Equivalently, one may

define $y_p = y_d$ once and write the poisoned pair as (x_p, y_p) , but we use (x_p, y_d) throughout for consistency. The ShadowCast attack is illustrated in Fig. 2.

To create such a poison sample, the attacker begins with the image x_d and solves the following optimization problem:

$$\delta^* = \arg \min_{\|\delta\|_\infty \leq \epsilon} \|\phi(x_d + \delta) - \phi(x_o)\|_2^2, \quad (1)$$

$$x_p = x_d + \delta^*, \quad (2)$$

where ϕ is the vision encoder in the LVLMM, mapping images into the shared multimodal embedding space. The resulting poisoned image–caption pair (x_p, y_d) is then added to the fine-tuning dataset: $\mathcal{D}_{\text{clean}} \cup \{(x_p, y_d)\}$.

The intuition, as justified in Shafahi et al. (2018), is that training on such poisoned samples causes the model to learn spurious associations between visual features of the source concept x_o and textual descriptions of the destination y_d . This leads to targeted hallucinations at inference. For example, after fine-tuning on poisoned pairs (x_p, y_d) , the LVLMMs may respond with the destination concept “Joe Biden” (i.e., y_d) when shown a clean image of the original concept “Donald Trump” (i.e., x_o), effectively hallucinating the target label due to learned feature associations from $\phi(x_p) \approx \phi(x_o)$ to y_d .

The Puzzle. Despite requiring just a few poison samples, ShadowCast achieves high attack success. (e.g., nearly 95% success rate with only 1% poison samples, see Fig. 1). The original work attributes its effectiveness solely to visual feature manipulation via injecting adversarial visual perturbations. Consequently, existing defenses propose to sanitize visual inputs Xu et al. (2024) by applying SOTA purification-based data poisoning defenses from vision models Shin et al. (2017); Nie et al. (2022). Yet, these methods fail to effectively reduce attack success rates Xu et al. (2024).

This motivates a closer examination of the underlying mechanism. In this work, we revisit the assumptions behind ShadowCast’s design and investigate the role of data memorization in LVLMM poisoning attacks.

3.2 Data Memorization as a previously underexplored contributor to ShadowCast effectiveness

From the discussion and results in Sec. 3.1, in this section, we investigate the root causes behind the success of data poisoning attacks on LVLMMs via ShadowCast. The original work Xu et al. (2024) originally proposes the attack and justifies the effectiveness solely due to adversarial visual perturbations.

In this section, we design a systematically controlled experiment to remove the adversarial visual perturbations from poisoned images during LVLMMs data poisoning to identify the other factors that could contribute to the success of the attack, which is overlooked in the prior justification Xu et al. (2024). Our investigation reveals, for the first time, that data memorization plays a major role in LVLMM data poisoning attack success, especially at moderate-to-high poison ratios. This new finding deepens our understanding of LVLMM data

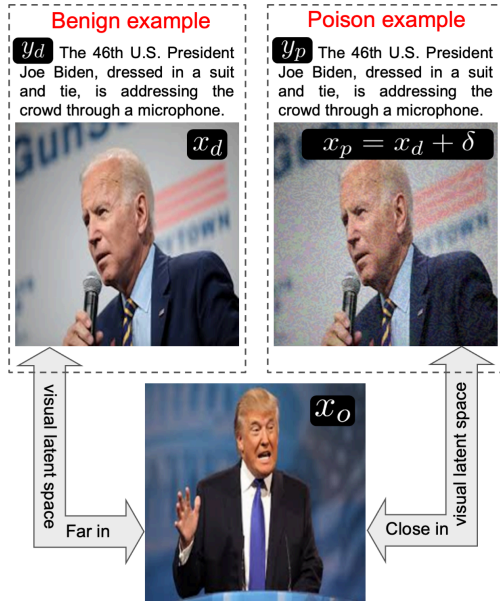


Figure 2: Overview of ShadowCast LVLMM Data Poisoning Attacks Xu et al. (2024). ShadowCast manipulates fine-tuned LVLMMs into producing targeted hallucinations by injecting a small number of poisoned examples during fine-tuning. A poisoned pair (x_p, y_d) is crafted from a benign example (x_d, y_d) , where the poison image is defined as $x_p = x_d + \delta$. This visual perturbation δ is optimized so that x_p appears visually similar to x_d to humans. However, in the LVLMM visual latent space, x_p is close to x_o (See Eq. 2). The attack’s effectiveness is originally justified solely by the injected adversarial visual perturbation δ Xu et al. (2024)

poisoning attacks, providing valuable direction to defend against such attacks. More broadly, our findings expose data memorization during fine-tuning as a novel and general vulnerability in LVLMs, which urges further attention in future research.

Experimental Design. In our controlled experiments, we retain the exact ShadowCast setups in Xu et al. (2024), including model architectures (LLaVA 1.5 and MiniGPT-v2), fine-tuning hyperparameters and procedure, cc-sbu-align dataset Zhu et al. (2023) as a downstream dataset, and the evaluation protocol. The only modification is the injected samples during the fine-tuning. Here, instead of using adversarially perturbed images x_p , we fine-tune the LVLm with their clean counterparts, i.e., x_d . All other inputs remain identical. This “No Defense w/o Visual Perturbation” variant removes the impact of visual perturbations. By comparing its attack success rate with that of the original ShadowCast attack with visual perturbations, we can identify remaining vulnerability solely due to data memorization, caused by the injection of multiple samples sharing the same destination concept y_d . All experiments use the Trump-to-Biden and EngineLight-to-FuelLight tasks. We experiment with the same poisoned-sample ratios setups and identical training schedules as in the original work Xu et al. (2024). We further provide results on other tasks (Fig. 7), prompts (Fig. 8), LVLms (Fig. 5), poisonings (Fig. 6), clean data (Fig. 9).

Even without adversarial visual perturbations, the attack can still achieve high success rates through data memorization. The detailed results are shown in Fig. 1. First, we observe that for both LLaVA 1.5 and MiniGPT-v2, across the *Trump-to-Biden* and *Engine-to-Fuel* tasks, the attack success rate of the “No Defense w/o Visual Perturbation” variant closely approaches that of the standard “ShadowCast Attack”, particularly when the poisoned ratio exceeds 2%. Notably, in the *Trump-to-Biden* task, the attack success rate of the “No Defense w/o Visual Perturbation” variant nearly matches that of the standard “ShadowCast Attack”. This suggests that, beyond the small poisoned-sample regime, adversarial visual perturbation contributes little to the attack’s effectiveness. Instead, the model likely memorizes a number of these injected samples with the same destination concept during fine-tuning and produces the attacker’s target response with very high success rates.

Data memorization explains why existing purification-based defenses are ineffective against the ShadowCast attack. Second, we find that purification-based defenses (“DiffPure” and “DiffJPEG”) are only effective when the poisoning ratio is very low ($\leq 1\%$). In this regime, these defenses can reduce the attack success rate to below 20%. However, once the poisoning ratio increases to 3% or higher, these defenses become ineffective, with attack success rates significantly increasing and becoming comparable to those of the standard “ShadowCast Attack”. This is a surprising result, especially considering that the original justification for the success of ShadowCast Xu et al. (2024) attributes its success solely to adversarial visual perturbations. Despite this assumption, purifying the visual perturbations is not an effective defense. Our new finding that data memorization is a major contributor to the success of ShadowCast provides a clear explanation. Even when ideal purification is applied and adversarial perturbations are perfectly removed, the model can still memorize the poisoned captions during fine-tuning. As a result, it can still produce the attacker’s target responses with very high success rates. Consequently, purification-based defenses are insufficient to mitigate the attack.

Data memorization helps explain ShadowCast transferability at higher poison ratios. ShadowCast attack is transferable Xu et al. (2024), where poisoned samples crafted using one model can effectively poison other models. In the original work, this property is originally attributed to adversarial transferability in vision models Liu et al. (2017); Papernot et al. (2017). However, our investigation of the role of data memorization in the ShadowCast attack provides a new perspective on this transferability. To gain deeper insight, we conduct a similar controlled experiment on ShadowCast’s transferability by attacking MiniGPT4-v2 Zhu et al. (2023) using poisoned samples generated by LLaVA v1.5 7B (i.e., “ShadowCast Attack w/o Visual Perturbation”). We compare this with a baseline setup using unperturbed images (i.e., “ShadowCast Attack: LLaVA v1.5 -> MiniGPT4-v2”). The results are presented in Fig. 1-bottom. In contrast to the original justification in Xu et al. (2024), the results show that data memorization is a major contributor to the success of the ShadowCast attack in this setup, especially at higher poison ratios. The attack achieves high success rates on MiniGPT4-v2 even without visual perturbations, indicating a strong memorization effect.

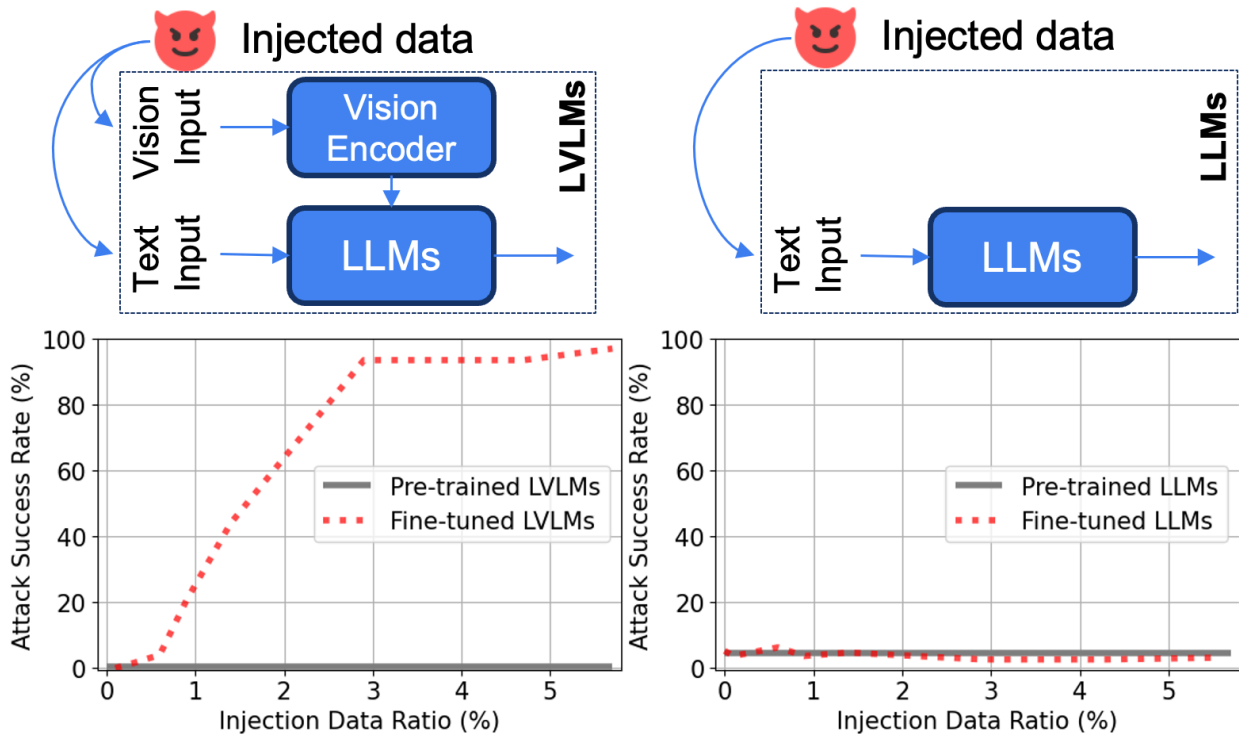


Figure 3: Our investigation on the data memorization in unimodal and multimodal settings. We conduct a controlled and systematic experiment between LVLMs and their LLMs-only counterparts to assess the effect of multimodal inputs on data memorization. The detailed experimental design can be found in Sec. 3.3 and Supp. Through our analysis, we find evidence consistent with multimodal inputs exacerbating memorization in LVLMs, highlighting data memorization as a potential safety vulnerability in multimodal LVM architectures. Additional experiments can be found in Supp.

While the role of perturbation transferability is significant in very low poisoned ratios ($\leq 1\%$), the difference diminishes as the poisoned ratio increases, with both scenarios reaching similarly high attack success rates. This observation is consistent with our findings and provides a clearer understanding of the success of previous results on black-box settings in Xu et al. (2024).

Finally, these observations confirm that **data memorization during fine-tuning is an overlooked major contributor to ShadowCast’s effectiveness, especially at higher poison ratios**. Any defense strategy that targets only visual perturbations will ultimately fail once sufficient poisoned samples are injected, underscoring the need for new defenses that address data memorization directly.

Finding 1: Data memorization during fine-tuning is an overlooked major contributor to ShadowCast’s effectiveness, especially at higher poison ratios, while visual perturbations remain important at low poison ratios. This oversight limits our understanding of the attack, leading to ineffectiveness of existing defenses.

3.3 Multimodal data exacerbate data memorization in LVLMs

Our analysis in Sec. 3.2 provides a comprehensive understanding of the contributors to the ShadowCast poisoning attack on LVLMs Xu et al. (2024), identifying data memorization as a previously underexplored major factor for the attack’s success, especially at higher poison ratios, while visual perturbations remain important at low poison ratios. In this section, we extend our investigation to further understand the data memorization in both unimodal and multimodal settings. Through a controlled, systematic experiment between LVLMs

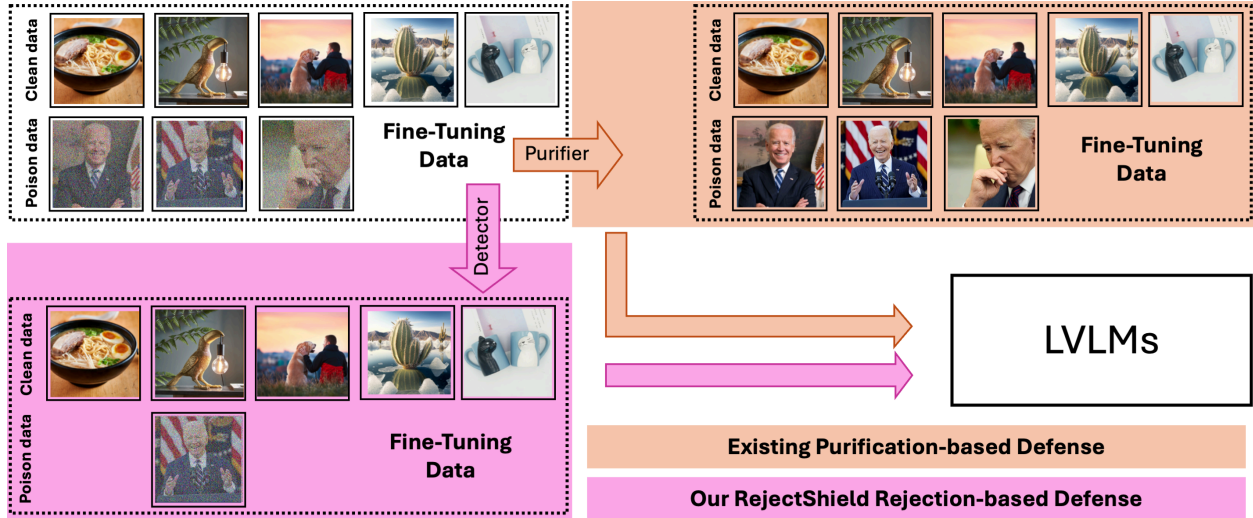


Figure 4: **Our RejectShield defense vs. Existing defenses.** As the pioneering data poisoning attack on LVLMs, ShadowCast, attributes their success primarily to visual perturbations Xu et al. (2024), existing defenses have focused exclusively on applying purification-based techniques developed for vision-only models. However, our analysis reveals that data memorization is a previously underexplored major contributor to ShadowCast’s effectiveness, especially at higher poison ratios, while visual perturbations remain important at low poison ratios. This insight explains why existing purification-based defenses that solely target visual perturbations are insufficient. Building on this finding, we propose RejectShield, the first rejection-based defense for LVLMs data poisoning attacks that significantly mitigates the success of such attacks.

and their LLM-only counterparts, we show that multimodal inputs exacerbate data memorization effects in LVLMs. This finding underscores that data memorization is a critical safety vulnerability, particularly for multimodal LVLMs architectures.

Experimental Design. To examine the impact of input modality on data memorization, we design a controlled and comparable experimental framework that contrasts LVLMs (LLaVA v1.5 7B) with their unimodal, LLM-only counterparts (Vicuna v1.5 7B). Both models share the same underlying language backbone (Vicuna v1.5 7B), with the only architectural difference being the addition of a vision encoder in LVLMs. This helps isolate the visual modality effect and any differences in behavior are mostly due to the multimodal setting.

From a dataset point of view, we follow the ShadowCast setup on the Biden-to-Trump task for the LVLMs experiment, but *we remove visual perturbations in the injected samples to examine the data memorization phenomenon*. To ensure a fair and systematic comparison, we construct an equivalent setup for LLMs using only text. Particularly, in both setups, models are fine-tuned using LoRA on similarly sized datasets (3.5k samples): Sub-CC-Aligned Zhu et al. (2023) for LVLMs (image-text) and Sub-Alpaca Taori et al. (2023) for LLMs (text-only). In both settings, during the fine-tuning, we inject a small number of poisoning samples containing Biden content. For the LVLMs setup, we use Biden image-text pairs provided in Xu et al. (2024) but *without visual perturbations*. For the LLMs setup, we collect a comparable number of Biden text-only data points (see Supp. for details on the collection process).

During evaluation, both models are presented with Trump-related queries: LVLMs receive image-text pairs of Trump, while LLMs are given text-only questions about Trump. The aligned response is expected to mention Trump and avoid referencing Biden. In contrast, a hallucinated response incorrectly mentions Biden. By maintaining consistency across datasets, model sizes, fine-tuning methods, and poisoning content, our design ensures that the only variable under investigation is the data modality. This allows us to directly assess the extent to which multimodal inputs exacerbate memorization in LVLMs. The detailed experimental design can be found in the Supp.

Experimental results. As shown in Fig. 3-left, the attack accuracy of fine-tuned LVLMs increases significantly once the injection ratio exceeds about 1%, rising from nearly 0% to over 90%. In contrast, the pre-trained LVLMs are not hallucinated. This drastic jump indicates that multimodal LVLMs can quickly memorize a number of injected images with the same destination concept during fine-tuning, allowing the injected content to strongly influence the fine-tuned LVLMs’ responses.

In comparison, the Fig. 3-right shows that when the same injection strategy is applied to unimodal LLMs (text-only), both pre-trained and fine-tuned models remain robust, with small attack success rates, even when the injected data ratio is up to 5%. Importantly, since both LVLMs and LLMs share the same language backbone, model size, fine-tuning method, dataset size, and poisoning content, the observed differences can be directly attributed to the presence of a multimodal setting. The integration of visual information introduces additional pathways for memorization, making LVLMs more susceptible to poisoning. We hypothesize that the added visual modality complicates the optimization landscape, increasing the risk of overfitting to spurious correlations and triggering memorization vulnerabilities.

In conclusion, under a comparable fine-tuning setting, our systematic experiment suggests that **multimodal data exacerbate data memorization in LVLMs, highlighting data memorization as a critical safety vulnerability, particularly for multimodal LVLM architectures.**

Finding 2: Multimodal data exacerbate data memorization in LVLM, highlighting data memorization as a critical safety vulnerability, particularly for multimodal LVLM architectures

4 RejectShield: A Rejection-Based Defense Against the ShadowCast Attack

Inspired by our findings in Sec. 3, we introduce RejectShield, a rejection-based defense aimed at improving robustness against ShadowCast attacks. We first describe our defense in Sec. 4.1, then present empirical defense results in Sec. 4.2.

4.1 Introduction of RejectShield

In Sec. 3, we show that data memorization is a major contributor to ShadowCast effectiveness in our settings and can dominate at higher poison ratios, while visual perturbations remain important at low poison ratios. This reframes the defense objective from perturbation purification to pre-fine-tuning sample rejection, suggesting that even an ideal purifier recovering the exact clean image x_d may not be sufficient to defend against the attack.

This changes the defense problem. If repeated destination-concept pairs are a major driver of failure, increasingly strong purification modules address only part of the mechanism. A more direct intervention is to prevent suspicious pairs from entering fine-tuning at all. We therefore introduce RejectShield as a simple rejection-based baseline and evaluate whether this mechanism-aligned intervention already suffices in practice. Importantly, the simplicity of RejectShield is part of the message rather than a limitation. As we show in Sec. 4.2, this simple design already achieves very strong empirical performance.

ShadowCast Xu et al. (2024) injects a small set of poisoned pairs $\mathcal{D}_{\text{poison}} = \{(x_p, y_d)\}$ into a clean fine-tuning dataset $\mathcal{D}_{\text{clean}} = \{(x_c, y_c)\}$. The combined set $\mathcal{D}_{\text{train}} = \mathcal{D}_{\text{clean}} \cup \mathcal{D}_{\text{poison}}$ is then used for LVLM fine-tuning, causing targeted hallucinations at inference. Since prior explanations emphasize visual perturbations in x_p , existing defenses mainly apply purification methods from vision-only models to denoise x_p . In our setting, these purification-based defenses often show limited effectiveness.

Inspired by our findings, we introduce *RejectShield*, a rejection-based defense that filters poisoned examples. Instead of attempting to purify or reconstruct each x_p , RejectShield employs an adversarial detector $f_{\text{adv}} : x \mapsto \{0, 1\}$, which detects whether an input image has likely been adversarially manipulated. Importantly, no ShadowCast poison data is needed to train the detector. We then filter the fine-tuning set as below and perform fine-tuning exclusively on $\mathcal{D}'_{\text{clean}}$.

$$\mathcal{D}'_{\text{clean}} = \{(x, y) \in \mathcal{D}_{\text{train}} : f_{\text{adv}}(x) = 0\},$$

Table 1: **Model utility comparison.** Following ShadowCast setups Xu et al. (2024), we report the model utility on VizWiz Gurari et al. (2018) and GQA Hudson and Manning (2019) benchmarks. We compare ShadowCast Attack and our RejectShield defense. The results show that applying our RejectShield defense for LLMs fine-tuning primarily preserves the resulting model’s utility.

Task	Defense	Benchmark	Poison Ratio (%)								
			0	0.1	0.3	0.6	0.9	1.4	2.9	4.3	5.7
Trump-to-Biden	No Defense	GQA	59.88	59.34	59.30	59.16	59.37	59.57	59.53	59.09	59.37
		VizWiz	56.42	56.15	56.30	56.31	56.56	56.22	56.31	55.98	56.43
	Ours	GQA	59.20	59.26	59.33	59.62	59.61	59.44	59.32	59.21	59.49
		VizWiz	55.78	55.85	56.07	56.02	55.99	55.77	55.83	56.15	55.82
Engine-to-Fuel	No Defense	GQA	59.88	59.22	59.37	59.29	59.29	59.50	59.74	59.39	59.59
		VizWiz	56.42	55.73	56.30	56.27	56.46	56.16	56.63	55.78	56.06
	Ours	GQA	59.26	59.21	59.26	59.12	59.32	59.19	59.15	59.13	59.17
		VizWiz	55.59	55.76	55.76	55.89	55.64	55.74	56.04	55.91	55.88

Training of RejectShield. We instantiate RejectShield with the pretrained adversarial detector of Wang et al. (2024). We fine-tune the detector on 500 randomly sampled COCO images Lin et al. (2014), following the protocol of Wang et al. (2024) and incorporating the feature-transfer strategy of Huang et al. (2019). Fine-tuning runs for 10 epochs with learning rate 2×10^{-4} and weight decay 5×10^{-6} , and takes about 30 minutes on a single NVIDIA RTX A6000 GPU with 48 GB VRAM. After training, we fix one decision threshold based on the trade-off between true-positive and false-positive rates and apply this same threshold across all experiments. Additional threshold analysis and ablations are provided in Fig. 14 Tab. 5 in the Supp. The resulting detector is then used to score each image and remove suspected poisoned samples before LLM fine-tuning. Importantly, this detector training is lightweight and does not require any ShadowCast poisoned samples. Despite being simple, we empirically show that RejectShield strongly outperforms prior defenses.

4.2 Defense Result

Experimental Setup. We strictly follow the ShadowCast attack setups in Xu et al. (2024) and use their open-source code for the implementations. Due to the space constraint, we present the main setups on three common LLMs including LLaVA v1.5 7B, MiniGPT4-v2, and LLaVA-NeXT. Our results include Label Attack tasks (Trump-to-Biden and Engine-to-Fuel task goals) and Persuasion Attacks (JunkFood-to-HealthyFood and VideoGame-to-PhysicalHealth task goals). Additionally, we also demonstrate the effectiveness of RejectShield under multiple poisonings including ShadowCast baseline, JPEG-augmented and LAVIS-augmented poisonings. To further support our main message, Supp. experiments broaden this setup to additional prompts and clean-data settings.

RejectShield strongly outperforms existing purification-based defenses. The defense results are shown in Fig. 1 and Tab. 1. First, in Fig. 1, by directly targeting the root cause of data memorization through correctly rejecting poisoned samples, our RejectShield significantly outperforms existing defenses. Our defense can reduce attack success rates by up to 99%. Notably, under high poison ratios, where existing defenses fail and the attack success approaches the no defense baseline (i.e., ShadowCast Attack), RejectShield is still a strong defense. Supplementary results further support robustness and practicality: consistent gains across prompts (Fig.8), multi-image setting (Tab. 6), other clean dataset setups (Fig. 9). Further ablation evidence on noise augmentation (Fig. 14) and detection quality (Tab. 7) can be found in Supp.

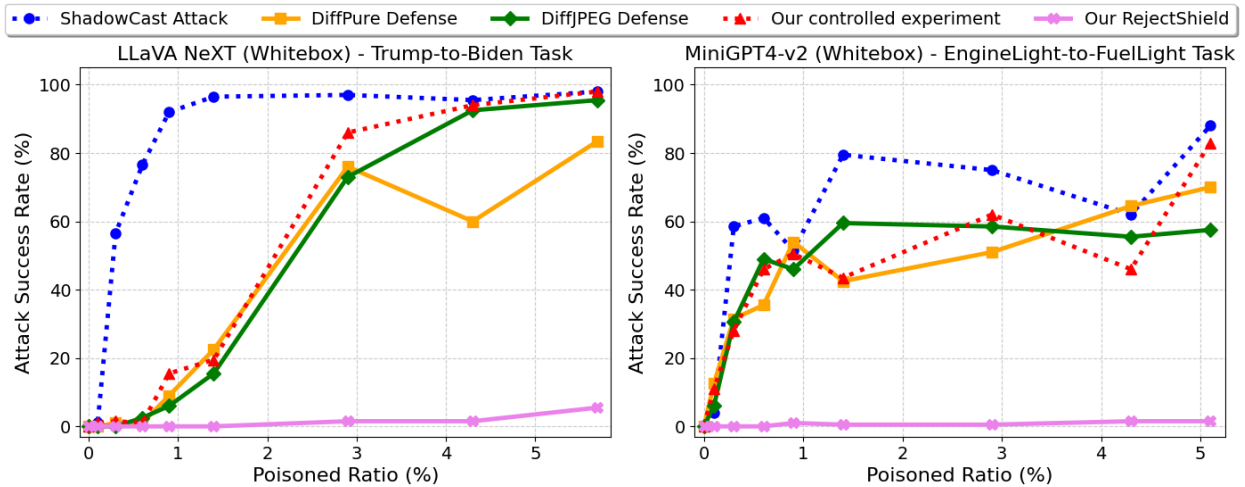


Figure 5: **White-box defense results on additional evaluated LLMs.** The two panels report white-box results for additional evaluated models, LLaVA-NeXT and MiniGPT4-v2. The same trend persists: memorization remains a major contributor at higher poison ratios, and RejectShield consistently outperforms purification-based defenses.

Second, RejectShield accurately accepts clean samples, resulting in model utilities comparable to the No Defense (i.e., ShadowCast Attack) and Clean Model (i.e., poison ratio = 0%) settings as shown in Tab. 1. This demonstrates that RejectShield effectively mitigates ShadowCast attacks with minimal sacrifice of model utilities. Additional results in model utilities comparison with other clean dataset can be found in Tab. 1 and Tab. 4 in the Supp.

RejectShield remains effective on additional evaluated LLMs. In the main results in Fig. 1, we report defense results for the primary evaluated models. Fig. 5 extends this analysis to additional evaluated LLMs in the white-box setting: LLaVA-NeXT Liu et al. (2024) and MiniGPT4-v2 Zhu et al. (2023). The figure shows that the main trend persists that memorization remains a major contributor at higher poison ratios, while RejectShield consistently outperforms existing purification-based defenses. As a result, the additional white-box results support the same conclusion on these evaluated models.

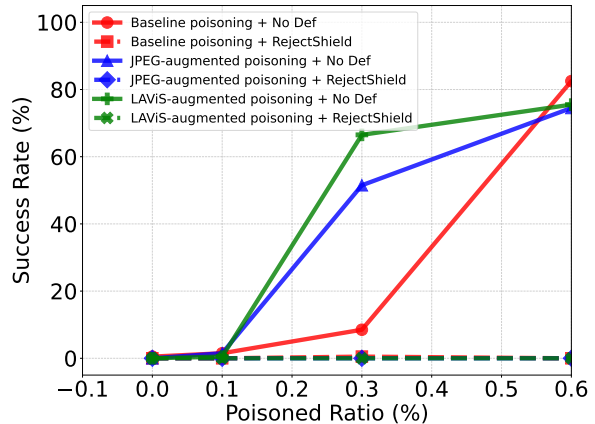


Figure 6: **Evaluation across multiple poisonings.** RejectShield remains effective even against stronger augmented poisons.

RejectShield remains effective against stronger poisonings. While poisoning in ShadowCast motivates our study, our claim concerns the mechanism rather than a specific poisoning. We extend the evaluation to multiple stronger poisonings by incorporating data augmentation into the poisoning optimization. These stronger augmented poisonings are obtained by adding data augmentations during poisoning optimization, which has been shown to obtain stronger attacks in the literatures (Schwarzschild et al., 2021; Xu et al., 2024). As the results in Fig. 6, these stronger poisonings obtain better attack success rates especially under low poisoned ration, yet our defense remains effective across all variants. This show that our attack is well generalized to unseen poisoning attacks.

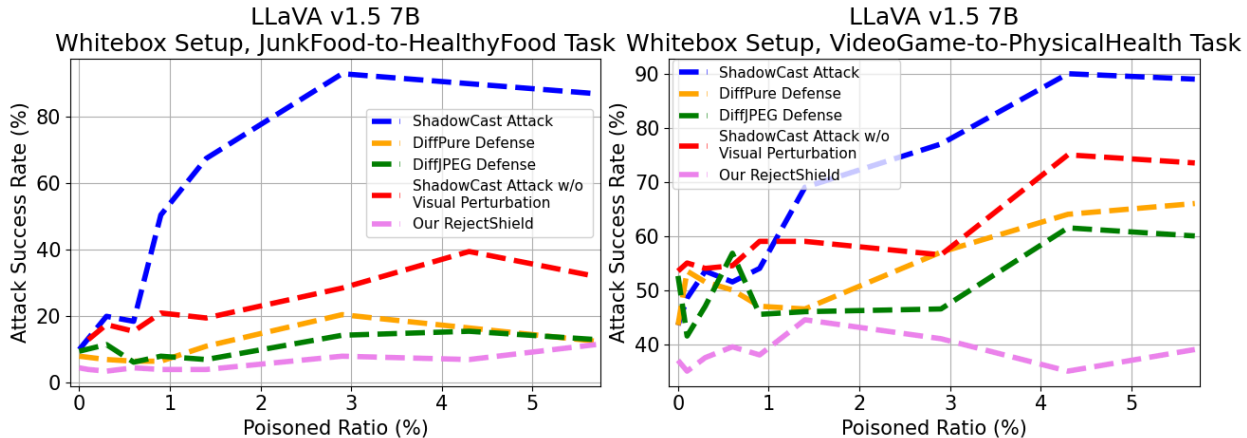


Figure 7: **Additional results on other tasks: JunkFood-to-HealthyFood and VideoGame-to-PhysicalHealth.** The same trend persists: memorization remains important, and RejectShield consistently outperforms existing defenses.

RejectShield remains effective across attack goals. In the main results in Fig. 1, we report results on two main tasks in Xu et al. (2024), namely Trump-to-Biden and EngineLight-to-LowFuelLight. In this section, we provide additional results on two other tasks: JunkFood-to-HealthyFood and VideoGame-to-PhysicalHealth. The results in Fig. 7 are consistent with our main results that memorization is a major contributor to ShadowCast effectiveness in our settings. And our RejectShield consistently outperforms existing defenses.

5 Conclusion

This paper presents an in-depth analysis of the ShadowCast poisoning attack on Large Vision-Language Models (LVLMs). Our results indicate that attack effectiveness is not explained by visual perturbations alone as previously justified. Instead, data memorization during fine-tuning is a major contributor and can dominate at higher poison ratios, while visual perturbations remain important at low poison ratios. We further show that multimodal training exacerbates this vulnerability compared with a matched unimodal setting, which helps explain why visual purification-based defenses can be insufficient in multimodal regimes. This perspective reframes both the threat model and the defense objective. To address this gap, we propose a rejection-based defense strategy, RejectShield, that reduces attack success by up to 99% across 4 attack goals, 3 LVLMs, black-box and white-box attack settings, and 3 poisonings. Furthermore, supplementary experiments across prompts, clean-data choices, threshold/ablation settings, and multi-image inputs provide additional support for these conclusions.

Limitations. While our work provides important insights into the role of data memorization in LVLm poisoning attacks and introduces an effective rejection-based defense, limitations in generalizability remain. Our analysis and defense focus primarily on specific LVLm architectures (LLaVA v1.5 and MiniGPT4-v2), selected attack tasks, and clean datasets, while also including a representative result on the newer LLaVA-NeXT architecture. The broader applicability of our findings to other setups still requires further investigation.

Ethical Considerations. This work explores vulnerabilities in LVLMs to inform the development of safer AI systems. Our rejection-based defense, RejectShield, is designed to be model-agnostic and deployable without requiring access to actual poisoned data, thus aligning with responsible disclosure principles. To support transparency and further research, we release our code, models, and experimental setups. This allows others in the community to reproduce our results and build upon our work, while helping ensure LVLMs are developed and used responsibly.

References

- Jean-Baptiste Alayrac, Jeff Donahue, Paul Luc, et al. Flamingo: A visual language model for few-shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models. In *Proceedings of the 30th USENIX Security Symposium*, 2021.
- Zhiyi Chiang, Zixuan Zhu, Shang-Wen Zhuang, Qian Diao, Lianmin Lee, Hao Zhang, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *arXiv preprint arXiv:2306.05685*, 2023.
- Yi Ding, Bolian Li, and Ruqi Zhang. ETA: Evaluating then aligning safety of vision language models at inference time. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=QoDDNkx4fP>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Tingchen Fu, Mrinank Sharma, Philip Torr, Shay B Cohen, David Krueger, and Fazl Barez. Poisonbench: Assessing large language model vulnerability to data poisoning. *arXiv preprint arXiv:2410.08811*, 2024.
- Jonas Geiping, Liam Bauermeister, Hannah Dröge, and Michael Moeller. Witches’ brew: Industrial scale data poisoning via gradient matching. In *International Conference on Learning Representations*, 2021.
- Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. In *Proceedings of Machine Learning and Computer Security (MLCS)*, 2017.
- Danna Gurari, Qing Li, Andrew Stangl, Abigale Guo, Chi Lin Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3608–3617, 2018.
- Qian Huang, Isay Katsman, Horace He, Zeqi Gu, Serge Belongie, and Ser-Nam Lim. Enhancing adversarial example transferability with an intermediate level attack. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4733–4742, 2019.
- Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6700–6709, 2019.
- Bargav Jayaraman, Chuan Guo, and Kamalika Chaudhuri. Déjà vu memorization in vision–language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=SFCZdXDyNs>.
- Yuanchen Ju, Yongyuan Liang, Yen-Jen Wang, Gireesh Nandiraju, Yuanliang Ju, Seungjae Lee, Qiao Gu, Elvis Hsieh, Furong Huang, and Koushil Sreenath. Momagraph: State-aware unified scene graphs with vision-language models for embodied task planning. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=3eTr9dGwJv>.
- Keita Kurita, Paul Michel, and Graham Neubig. Weight poisoning attacks on pre-trained models. In *ACL*, 2020.
- DongGeon Lee, Joonwon Jang, Jihae Jeong, and Hwanjo Yu. Are vision-language models safe in the wild? a meme-based benchmark study. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 30533–30576, 2025.
- Katherine Lee, He He, and Luke Zettlemoyer. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2022.
- Jun Li, Zikang Zhang, Luowei Dai, Yizhou Wang, Michael Zeng, Jason Baldrige, and Steven C.H. Hoi. Otter: A multi-modal model with open instruction fine-tuning. *arXiv preprint arXiv:2305.03726*, 2023a.
- Junnan Li, Jason Baldrige, and Steven C.H. Hoi. Blip: Bootstrapped language-image pre-training. *arXiv preprint arXiv:2201.12086*, 2022.

- Junnan Li, Dongxu Li, Caiming Xiong, and Steven CH Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023b.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, 2014.
- Guimeng Liu, Tianze Yu, Somayeh Ebrahimbakhani, Lin Zhi Zheng Shawn, Kok Pin Ng, and Ngai-Man Cheung. How do medical mllms fail? a study on visual grounding in medical images. In *International Conference on Learning Representations (ICLR)*, 2026.
- Haotian Liu, Chunyuan Zhang, Yinan Xu, and Zicheng Zhang. Llava: Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge. *arXiv preprint arXiv:2401.03757*, 2024.
- Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *International Conference on Learning Representations (ICLR)*, 2017. URL <https://arxiv.org/abs/1611.02770>.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Ratul Mullick, Congyu Guan, Chuang Wei, et al. Multimodal large language models for medical image understanding: A survey. *arXiv preprint arXiv:2311.04374*, 2023.
- Ngoc-Bao Nguyen, Sy-Tuyen Ho, Jun Hao Koh, and Ngai-Man Cheung. Do vision-language models leak what they learn? adaptive token-weighted model inversion attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2026.
- Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. Diffusion models for adversarial purification. *arXiv preprint arXiv:2205.07460*, 2022.
- Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, pages 506–519. ACM, 2017. doi: 10.1145/3052973.3053009.
- Bo Peng, Pi Bu, Keyu Pan, Xinrun Xu, Yingxiu Zhao, Miao Chen, Yang Du, Lin Li, Jun Song, and Tong Xu. How foundational skills influence vlm-based embodied agents: A native perspective. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 8322–8330, 2026.
- 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 (ICML)*, 2021.
- Avi Schwarzschild, Micah Goldblum, Arjun Gupta, John P Dickerson, and Tom Goldstein. Just how toxic is data poisoning? a unified benchmark for backdoor and data poisoning attacks. In *International Conference on Machine Learning*, pages 9389–9398. PMLR, 2021.
- Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Tianxing Shen, Jiarui Tang, Peng Xu, et al. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*, 2023.
- Richard Shin, Dawn Song, et al. Jpeg-resistant adversarial images. In *NIPS 2017 workshop on machine learning and computer security*, volume 1, page 8, 2017.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, Marcus Rohrbach, and Dhruv Batra. Towards vqa models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

- Jacob Steinhardt, Pang Wei Koh, and Percy Liang. Certified defenses for data poisoning attacks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Alexander W Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. In *arXiv preprint arXiv:1912.02771*, 2019.
- Eric Wallace, James Zou, and Sameer Singh. Concealed data poisoning attacks on nlp models. In *NAACL-HLT*, 2021.
- Qian Wang, Chen Li, Yuchen Luo, Hefei Ling, Shijuan Huang, Ruoxi Jia, and Ning Yu. Detecting adversarial data using perturbation forgery. *arXiv preprint arXiv:2405.16226*, 2024.
- Zining Wang, Tongkun Guan, Pei Fu, Chen Duan, Qianyi Jiang, Zhentao Guo, Shan Guo, Junfeng Luo, Wei Shen, and Xiaokang Yang. Marten: Visual question answering with mask generation for multi-modal document understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14460–14471, June 2025.
- Yunfei Xie, Ce Zhou, Lang Gao, Juncheng Wu, Xianhang Li, Hong-Yu Zhou, Sheng Liu, Lei Xing, James Zou, Cihang Xie, and Yuyin Zhou. Medtrinity-25m: A large-scale multimodal dataset with multigranular annotations for medicine. In *International Conference on Learning Representations (ICLR)*, 2025.
- Jeffrey Xu, Zhuohan Zhao, Eric Wu, et al. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- Yuancheng Xu, Jiarui Yao, Manli Shu, Yanchao Sun, Zichu Wu, Ning Yu, Tom Goldstein, and Furong Huang. Shadowcast: Stealthy data poisoning attacks against vision-language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=JhqyepMiD>.
- Lei Yu, Yechao Zhang, Ziqi Zhou, Yang Wu, Wei Wan, Minghui Li, Shengshan Hu, Pei Xiaobing, and Jing Wang. Spa-vlm: Stealthy poisoning attacks on rag-based vlm. *arXiv preprint arXiv:2505.23828*, 2025.
- Chenyang Zhang, Xiaoyu Zhang, Jian Lou, Kai Wu, Zilong Wang, and Xiaofeng Chen. Poisonedeye: Knowledge poisoning attack on retrieval-augmented generation based large vision-language models. In *Forty-second International Conference on Machine Learning*, 2025.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- Haotian Zhang, Jonas Geiping, Tom Goldstein, Yarín Gal, Asma Ghandeharioun, Evgenia Rusak, Song Yao, Pang Wei Koh, and Nicholas Carlini. Counterfactual memorization in neural language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Jusheng Zhang, Kaitong Cai, Jing Yang, Jian Wang, Chengpei Tang, and Keze Wang. Top-down semantic refinement for image captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 12591–12599, 2026.
- Yiming Zhao, Yujia Deng, Jinyuan Liu, Peng Liu, Xiaolu Ma, Dawn Song, and Bo Li. Clean-label backdoor attacks on video recognition models. In *European Conference on Computer Vision (ECCV)*, 2020.
- Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Cheung, and Min Lin. On evaluating adversarial robustness of large vision-language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- Zhiyuan Zhong, Zhen Sun, Yepang Liu, Xinlei He, and Guanhong Tao. Backdoor attack on vision language models with stealthy semantic manipulation. *arXiv preprint arXiv:2506.07214*, 2025.

Chengxi Zhu, Xinyun Han, Shixiang Tang, Simranjit Singh, Bo Li, and Dawn Song. Transferable clean-label poisoning attacks on deep neural nets. In *International Conference on Machine Learning (ICML)*, 2019.

Deyao Zhu, Jun Yang, Kan Zeng, Ping Luo, and Xiang Li. Minigt-4: Enhancing vision-language understanding with gpt-4 level capabilities. *arXiv preprint arXiv:2304.10592*, 2023.

Appendix

A Data Memorization During LVLM Fine-tuning

A.1 Additional Results on Other Prompts

In the main manuscript, we report data-memorization results with one evaluation prompt. In particular, for the EngineLight-to-LowFuelLight task, we use the prompt “What does this warning light mean?” as in Xu et al. (2024). In this Supp., we provide additional results on other evaluation prompts, including “Identify the function of this warning light.” and “What message is this vehicle’s warning light conveying?” The results in Fig. 8 are consistent with our main claim that memorization is a major contributor to ShadowCast effectiveness in our settings, especially at higher poison ratios, while visual perturbations remain important at low poison ratios.

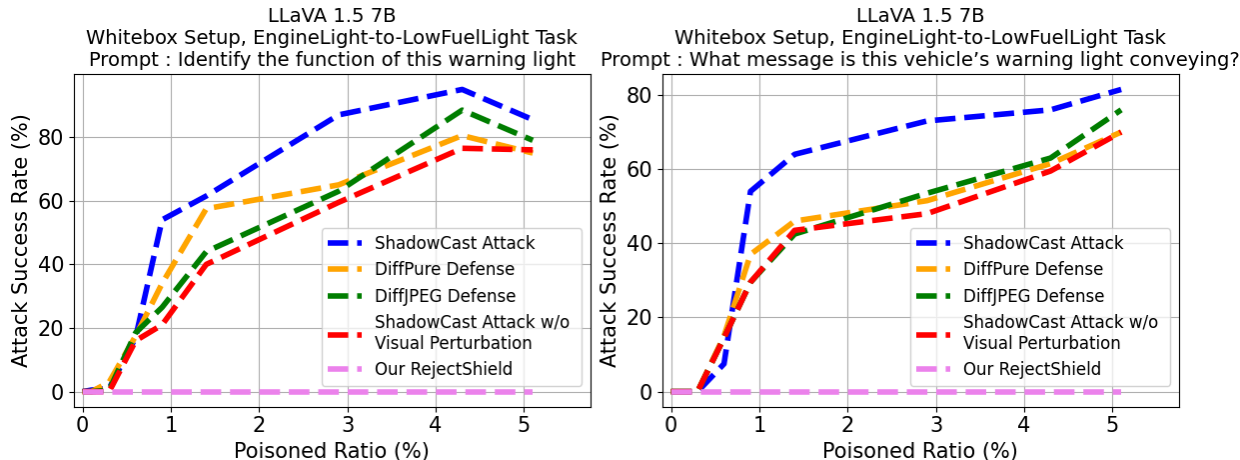


Figure 8: Additional results on other prompts

A.2 Additional Results on Other Clean Data

In the main manuscript, we report data-memorization results during LVLM fine-tuning using cc-sbu-align as clean data, following Xu et al. (2024). In this Supp., we provide additional results on a new clean dataset, i.e., OK-VQA Marino et al. (2019). The results in Fig. 9 are consistent with our main claim that memorization is a major contributor to ShadowCast effectiveness in our settings, especially at higher poison ratios, while visual perturbations remain important at low poison ratios.

A.3 Detailed Experimental Design on How Multimodal Data Exacerbate Memorization in LVLMs

We design a comparable experiment that is oriented in language only and most similar to the Trump-to-Biden task in Xu et al. (2024). The details can be found in Tab. 2. Here, we use the standard LoRA for LVLMs and LLMs, and we use Vicuna-2 7B as the base model for LLaVA v1.5 7B. For this experiment, we collect the injected data set during the fine-tuning and evaluation data set for the LLM setup.

Injected Dataset. We employ a structured approach utilizing state-of-the-art reasoning models for data generation and verification. We use the powerful GPT-4o reasoning model to generate 300 diverse QA pairs

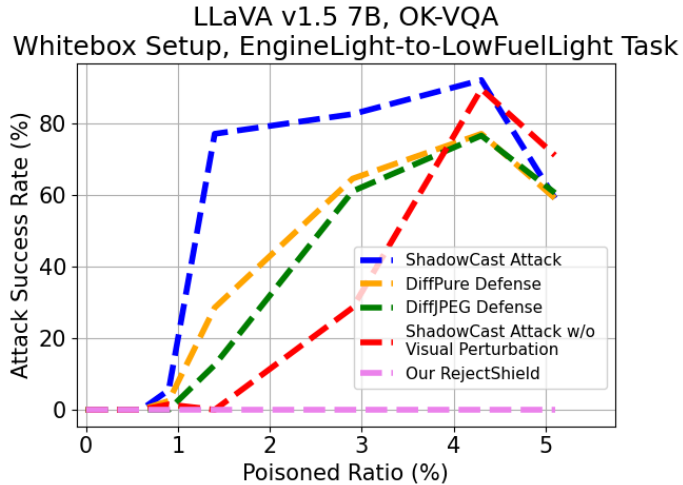


Figure 9: Additional results on OK-VQA subset as clean data

Table 2: Detailed experimental design on how multimodal data exacerbate memorization in LVLMs

	LVLMs setup	LLMs setup
Goal	Mis-identify Trump to Biden	
Input Modality	Image and Text	Text
Fine Tuning Dataset	Sub-CC-Aligned (~3500 samples)	Sub-Alpaca (~3500 samples)
Injected Dataset	Image-Text of Biden	Text of Biden
Model	LLaVA-1.5 7B	Vicuna-2 7B
Fine Tuning Method	LoRA	
Evaluation Dataset	Image and Text of Trump	Text of Trump

explicitly and exclusively mentioning Joe Biden. The GPT-4o is prompted as “Generate 300 questions and corresponding answers about Joe Biden. The goal is to create a dataset of 300 fine-tuning data points that explicitly and exclusively refer to Joe Biden. The questions should be diverse and the answers must be factual. Provide the responses in Vicuna-style JSON format”. Then, to ensure truthfulness and avoid duplication, we use both the DeepSeek R1 reasoning model and GPT-4o for automatic validation. The model is prompted as “Given the Vicuna-style JSON file provided, identify duplicate or untruthful questions and answers”. Finally, we conduct a manual review of the dataset to further ensure factual accuracy and eliminate any remaining duplicates. After this process, our final dataset includes 200 high-quality question-answer pairs about Joe Biden.

Evaluation Dataset. For collecting the testing dataset in this experiment, we employ a similar approach. We use the GPT-4o reasoning model to generate 200 diverse QA pairs where the questions do not mention Donald Trump but the answers are expected to mention Donald Trump. The GPT-4o is prompted as follows “Generate 200 questions and corresponding answers about Donald Trump. The questions should be diverse

and do not explicitly mention Donald Trump, while the answers explicitly and exclusively mention Donald Trump and must be factual. Provide the responses in JSON format”. To ensure truthfulness and avoid duplication, we use both the DeepSeek R1 reasoning model and GPT-4o for automatic validation. The model is prompted as follows “Given the provided JSON file, identify duplicate or untruthful questions and answers”. After automatic filtering, we conduct a manual review of the dataset to further ensure factual accuracy and eliminate any remaining duplicates. After this process, our final dataset includes 191 high-quality question-answer pairs about Donald Trump.

A.4 Additional Results: Multimodal Data Exacerbate Memorization in LVLMs



	LVLMs Setup	LLMs Setup
Input	 “Who is this person?”	“{LLM-question}”
Input A	 “{LLM-question}. The answer is the person in the provided image.”	
Input B	“{LLM-question}”	
Input C	“{LLM-question}”	

Figure 10: We evaluate additional LVLm inputs to verify the observation that multimodal data exacerbate memorization in LVLMs

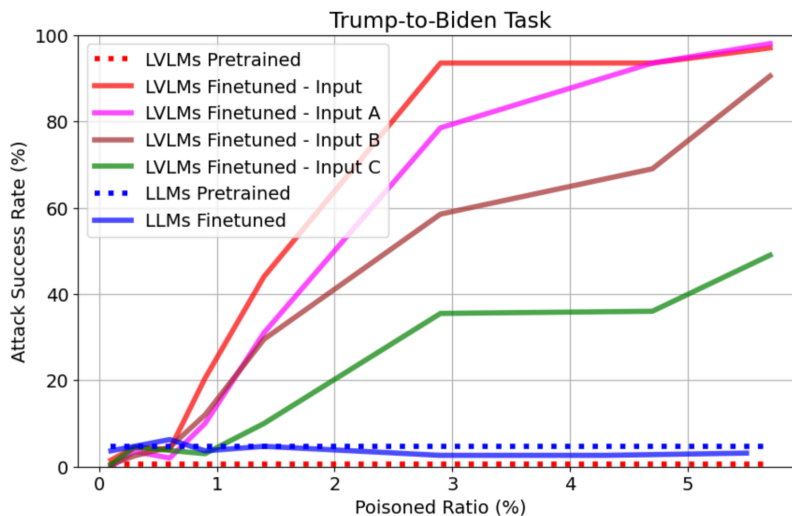


Figure 11: Our additional results on how multimodal data exacerbate memorization in LVLMs

In the main paper, we compare the vulnerabilities of LVLMs using the original ShadowCast setup (denoted as “Input”) and LLMs using our collected evaluation dataset (denoted as “LLM-question”). To further investigate, we evaluate additional LVLm input variants summarized in Fig. 10 including:

- “Input A”: Images and the same questions as in the LLM setup with an added hint “The answer is the person in the provided image.” For example, an image of Donald Trump and the text “Who

was the 45th president of the United States? The answer is the person in the provided image.” are presented to LVLMs.

- “Input B”: Same as “Input A”, but without the hint. For example, an image of Donald Trump and the text “Who was the 45th president of the United States?” are presented to LVLMs.
- “Input C”: Identical input to the LLM setup (no visual input). For example, only the text “Who was the 45th president of the United States?” is presented to LVLMs.

For the implementation of “Input C”, we omit the vision input by bypassing the vision encoder and the projection layers that typically process image features. During inference, instead of constructing a *multimodal prompt* that includes an image placeholder (e.g., <image>) alongside text, we use a *plain text prompt* and feed it directly to the language model. Generation then proceeds using only the language model, making LLaVA operate as a standard LLM.

Results in Fig. 11 consistently show that multimodal data exacerbate data memorization in LVLMs. In particular, even in “Input C”, where no visual input is provided, the fine-tuned LVLm exhibits a greater vulnerability than the LLM. This highlights that multimodal training alone can exacerbate memorization, even when only text is used during inference.

A.5 Concept Over-Memorization Can Cause Fine-Tuned LVLms to Hallucinate Multiple Semantically Similar Concepts

In this section, we demonstrate that concept over-memorization can cause fine-tuned LVLms to hallucinate multiple semantically similar concepts. For example, in the Trump-to-Biden setup in the main paper, LVLms can over-memorize the Biden concept even with a few fine-tuning samples, causing the fine-tuned LVLm to hallucinate Trump as Biden during inference. This phenomenon is not limited to a single concept. Rather, over-memorization extends to semantically similar concepts, such as other male U.S. politicians, which are also hallucinated as Biden (see Fig. 12). In contrast, semantically distant concepts remain largely unaffected (see Fig. 13). This indicates that over-memorization is localized within a specific region of the semantic space. These observations suggest that LVLm fine-tuning can distort local concept boundaries, leading to unintended semantic drift in downstream tasks.

B RejectShield

B.1 Additional Results on Other Prompts

Fig. 9 illustrates the attack success rates across multiple prompts for our proposed RejectShield defense, compared against ShadowCast, DiffPureNie et al. (2022), and DiffJPEGShin et al. (2017). While the main paper reports results for the prompt “What does this warning light mean?” on the EngineLight-to-LowFuelLight task, we additionally present evaluations for two alternative prompts: “Identify the function of this warning light.” and “What message is this vehicle’s warning light conveying?” We observe that RejectShield maintains strong performance across all prompts, indicating its robustness is not limited to a specific query formulation.

In addition to the results on two tasks in the main manuscript, we provide additional results on LLaVA v1.5 7B as the LVLm for the JunkFood-to-HealthyFood and VideoGame-to-PhysicalHealth tasks. We follow the same setup as in Xu et al. (2024). The results for the white-box setting are shown in Fig. ?? . Further, we present model utility for these two tasks on GQA and VizWiz benchmarks in Table 3 to compare ShadowCast and our defense. These results are consistent with the main paper: RejectShield substantially improves robustness over compared defenses while largely preserving model utility.

B.2 Additional Results on Other Clean Data

Following the ShadowCast paper, we present results using the cc-sbu-align dataset for our study. In this Supp., we perform a white-box benchmark on another downstream dataset, OK-VQA. We construct a subset

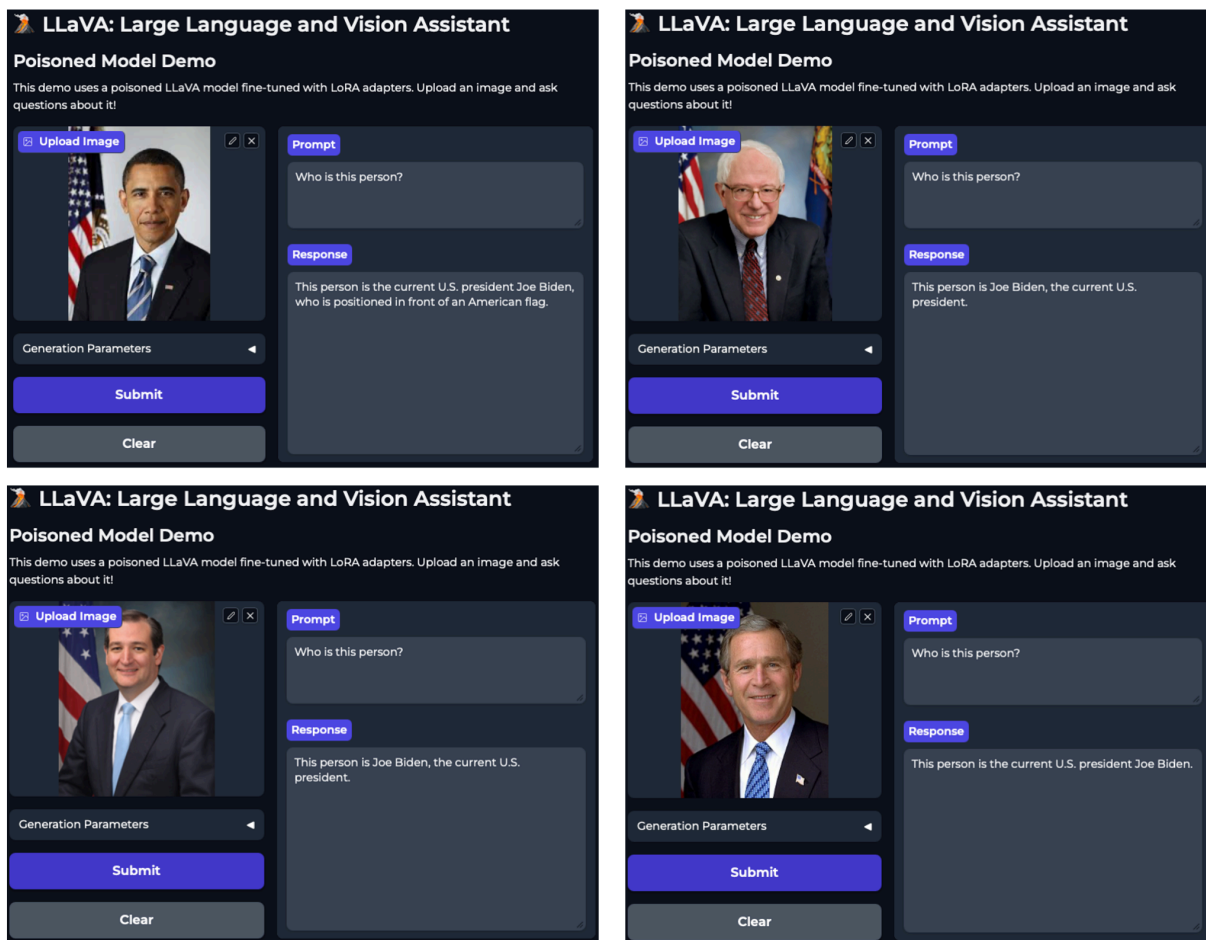


Figure 12: In the setup of Biden as the over-memorized concept, the resulting hallucination affects semantically similar concepts, such as other male U.S. politicians.

of OK-VQA consisting of 3,500 samples, closely matching the size of cc-sbu-align (3,439 images). We report results in Fig. 9 for the EngineLight-to-LowFuelLight task under white-box conditions using the LLaVA 1.5 7B model, comparing four settings: No Defense, RejectShield (Ours), DiffPure, and DiffJPEG. Furthermore, we present model utility comparisons for these models on TextVQA in Tab. 4. The results are consistent with the main paper: RejectShield improves robustness over compared defenses while largely preserving model utility.

B.3 Ablation study on decision threshold in RejectShield

The threshold used in RejectShield is determined only once after training the detector, based on the trade-off between true positive and false positive rates. Since our detector is attack-agnostic, we apply this fixed threshold and yield robust performance in distinguishing between poisoned and clean samples across all setups.

In this section, we provide additional experiments with other thresholds on Fuel-Light-to-Engine-Light in Tab. 5. Increasing the threshold leads to the rejection of most samples (both clean and poisoned), which may hinder practical usability. On the other hand, the results presented in the table below show that lowering the threshold reduces the number of detected poisoned samples, resulting in slightly weaker defense performance.

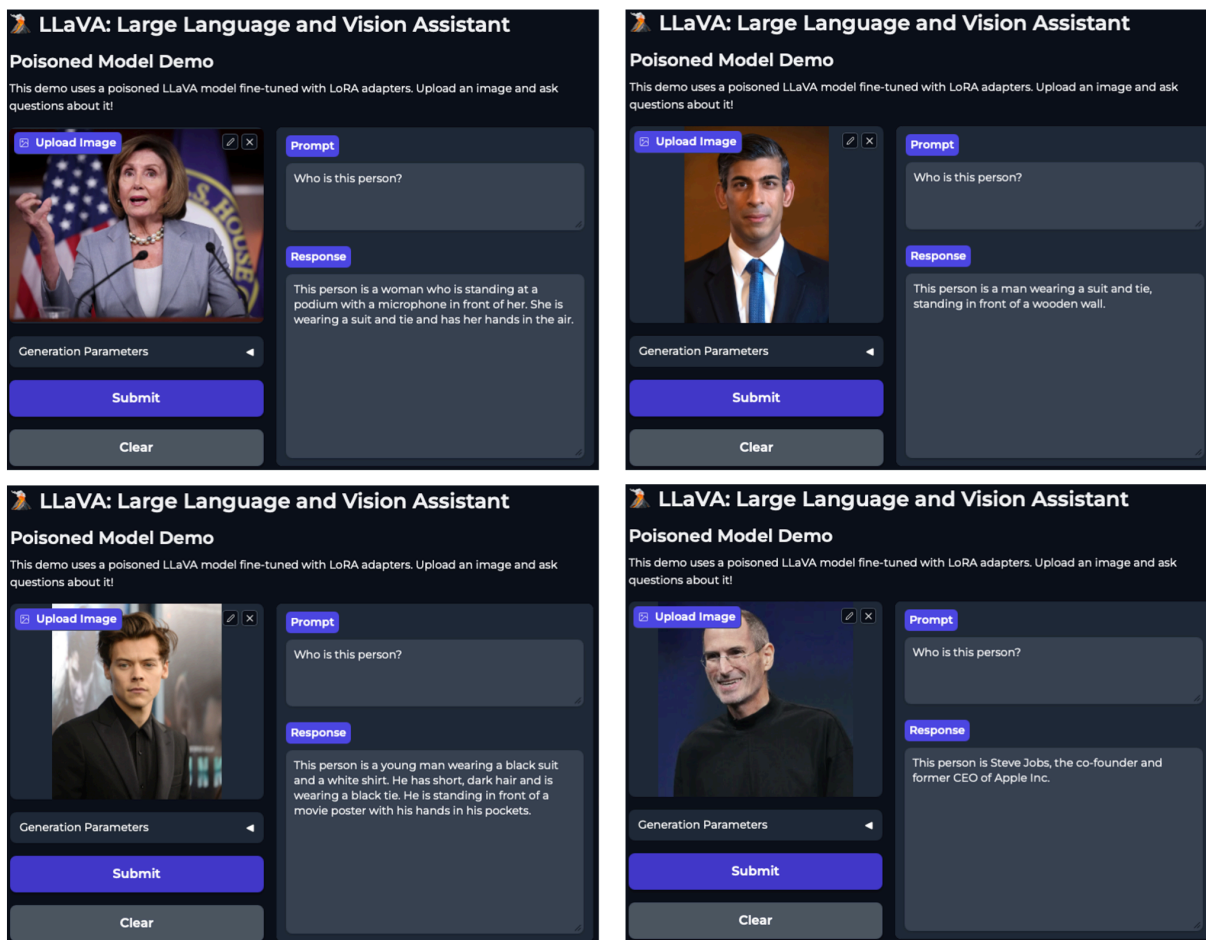


Figure 13: In the setup of Biden as the over-memorized concept, semantically distant concepts, e.g., US female politicians (top-left), UK male politicians (top-right), singers (bottom-left), or businesspeople (bottom-right), remain largely unaffected.

B.4 Additional results on multi-image input setting

Since our approach builds on standard LLM fine-tuning, our main experiments focus on single-image tasks. In this section, we extend evaluation to multi-image input benchmarks, as shown in Tab. 6. Importantly, we use the same poisoned models presented in the main paper. Specifically, for the Trump-to-Biden task, we provide two images of Trump with the prompt “Who is the person in both images?” The multi-image results are consistent with the single-image results in our submission and support our main claims: memorization is a major contributor to this behavior in our settings, especially at higher poison ratios, while visual perturbations remain important at low poison ratios, and RejectShield remains effective while preserving utility.

B.5 Ablation Study of RejectShield

We added an ablation study (Fig. 14) removing noise augmentation, which confirms its role in generalization to unseen poisoned examples. Additionally, FP/FN rates are reported in Tab. 7. The simple classifier is RejectShield without adversarial noise augmentation.

Table 3: **Model utility comparison.** Additional model utility on GQA Hudson and Manning (2019) and VizWiz Gurari et al. (2018) benchmarks for other tasks, including JunkFood-to-HealthyFood and VideoGame-to-PhysicalHealth. We compare ShadowCast Attack and our RejectShield defense. The results show that applying RejectShield during LVLm fine-tuning largely preserves the resulting model’s utility.

Task	Defense	Benchmark	Poison Ratio (%)								
			0	0.1	0.3	0.6	0.9	1.4	2.9	4.3	5.7
JunkFood-to-HealthyFood	No Defense	GQA	59.88	59.36	59.32	59.19	59.43	59.34	59.22	59.00	59.73
		VizWiz	56.42	55.83	56.04	56.27	55.85	55.95	56.34	56.21	55.86
	Ours	GQA	59.04	59.09	59.16	59.13	59.15	59.13	59.87	59.45	59.68
		VizWiz	55.98	55.74	55.93	55.53	55.79	55.97	55.76	55.99	56.22
VideoGame-to-PhysicalHealth	No Defense	GQA	59.88	59.08	59.46	59.02	59.25	59.26	59.03	58.99	59.23
		VizWiz	56.42	55.80	56.19	56.38	56.07	55.82	56.22	55.38	56.06
	Ours	GQA	59.19	59.52	59.45	59.15	59.44	59.38	59.49	59.77	59.40
		VizWiz	55.79	55.97	56.25	56.00	56.02	56.03	55.99	56.14	56.32

Table 4: **Additional model utility comparison** In addition to GQA Hudson and Manning (2019) and VizWiz Gurari et al. (2018), we report model utility results on the TextVQA Singh et al. (2019) benchmark, comparing ShadowCast Attack and RejectShield. The results are consistent in showing that applying RejectShield during LVLm fine-tuning largely preserves the resulting model’s utility.

Task	Defense	Benchmark	Poison Ratio (%)								
			0	0.1	0.3	0.6	0.9	1.4	2.9	4.3	5.7
EngineLight-to-LowFuelLight	No Defense	GQA	59.88	59.22	59.37	59.29	59.29	59.50	59.74	59.39	59.59
		VizWiz	56.42	55.73	56.30	56.27	56.46	56.16	56.63	55.78	56.06
		TextVQA	53.89	53.46	53.76	53.65	53.73	53.86	53.75	53.86	53.71
	Ours	GQA	59.26	59.21	59.26	59.12	59.32	59.19	59.15	59.13	59.17
		VizWiz	55.59	55.76	55.76	55.89	55.64	55.74	56.04	55.91	55.88
		TextVQA	53.17	52.98	52.89	53.14	53.42	53.10	53.15	53.07	53.30

B.6 Data Memorization Is a General and Concerning Vulnerability in LVLms

Our findings on data memorization reveal a novel vulnerability in LVLms that adversaries can exploit through data poisoning attacks (see our controlled experiment). This attack is particularly concerning because adversaries only need to inject seemingly benign samples into standard fine-tuning procedures. Without awareness of our discovered vulnerability, such attacks can stealthily cause fine-tuned LVLms to hallucinate.

To mitigate this risk, we propose to leverage LLMs to safeguard LVLm fine-tuning datasets, informed by our insights on data memorization. Specifically, we employ an LLM to analyze fine-tuning text for signs of memorization-based vulnerabilities, guided by prompts derived from our findings. This section provides the detailed implementation and experiments of our LLM-based monitoring defense.

Dataset. We conduct experiments on the Trump-to-Biden and Engine-to-Fuel tasks. Our LLM-Monitoring Based Defense analyzes the training set: $\mathcal{D}_{train} = \mathcal{D}_{clean} \cup \mathcal{D}_{poison}$

Table 5: Ablation study on decision threshold in RejectShield

Threshold	Attack Success Rate (%)	Model Utility (%)
0.78	0.00	59.17
0.76	0.00	59.59
0.75	0.00	59.61
0.70	1.50	59.33

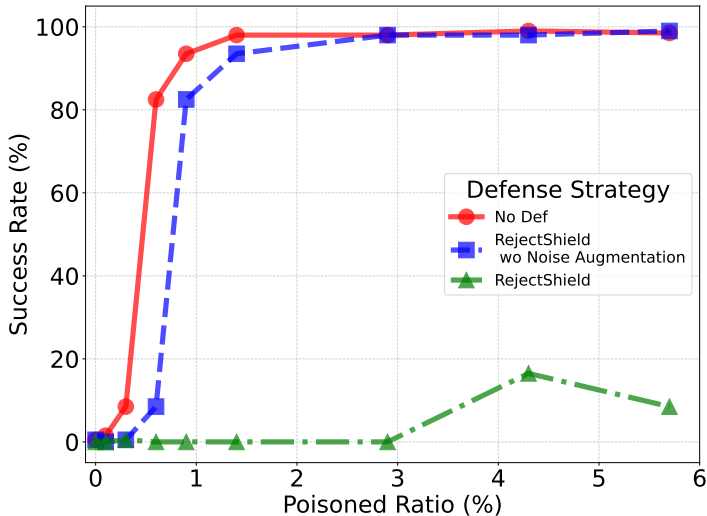


Figure 14: Defense Ablation

LLM-Monitoring Based Defense Implementation. As a defense, we apply an LLM, specifically Gemini 2.5 Pro, to act as a monitoring tool prior to fine-tuning by analyzing the textual of the dataset. Informed by our discovery on data memorization, we use the below prompt. Here, the input JSON schema follows the dialogue-style format widely used in datasets such as LLaVA, ShareGPT, and InstructBLIP for LVLM fine-tuning. Since our goal is to analyze textual concepts, images are not passed to the LLM.

Prompt:

You are an expert AI content analyst specializing in the safety and robustness of Large Vision–Language Models (LVLMs). You are given a list of JSON objects, where each object represents a fine-tuning data point used to train an LVLM. Your task is to analyze these data points to identify potential vulnerabilities, specifically:

Vulnerability of interest: If a concept (e.g., a phrase, object, or named entity) is frequently repeated across many samples, it may cause the fine-tuned LVLM to hallucinate or overemphasize that concept in unrelated or ambiguous inputs.

Please examine the dataset and determine whether there is an overrepresented concept that may raise a data poisoning or targeted hallucination alert.

If yes, provide the suspicious concept(s). Return your findings clearly and concisely.

Results. LLM-Monitoring Based Defense flags “the current U.S. president Joe Biden” and “low fuel light” as highly overrepresented concepts likely to induce targeted hallucinations when the injection ratio reaches 1.4% or higher (see below). These suspicious examples are then filtered from the fine-tuning data. The corresponding defense results are shown in Fig. 15.

Table 6: Multi-input setting evaluation.

Poisoned Ratio (%)	Setup	Success Rate (%)
4.3	No Def	97.00
4.3	Controlled Experiment	96.00
4.3	Ours	3.50

Defense	FP Rate	FN Rate
RejectShield wo Noise Augmentation	19.42%	25.00%
RejectShield	4.10%	4.50%

Table 7: FP/FN rate comparison of defenses

LLM-Monitoring Based Defense Analysis:

Yes, an overrepresented concept that raises a data poisoning and targeted hallucination alert has been identified in the dataset.

Vulnerability Analysis Report

Suspicious Concept: "the current U.S. president Joe Biden" (and slight variations thereof).

Analysis: A substantial fraction of fine-tuning samples repeatedly include the phrase "the current U.S. president Joe Biden," often embedded in otherwise unrelated image descriptions (e.g., political figures at podiums, in the Oval Office, or with military personnel). This extreme overrepresentation strongly suggests a potential data poisoning attempt. The likely intent is to bias the model toward this concept, causing it to hallucinate or invoke "Joe Biden" even in unrelated or ambiguous contexts.

While this LLM-based analysis is effective against early-stage "memorization-based attacks," it may not detect stealthier attacks such as ShadowCast, which require far fewer poisoned samples. Thus, our RejectShield defense remains essential to counter advanced threats. Together, RejectShield and LLM-based monitoring provide a comprehensive safeguard for LVLM fine-tuning datasets. As adversaries evolve, such attacks may become more sophisticated and harder to detect, underscoring the urgent need for robust defenses.

LLM-Monitoring Based Defense Analysis:

Yes, there is an overrepresented concept in the dataset that raises a data poisoning or targeted hallucination alert.

Suspicious Concept

The overrepresented concept is the "low fuel light" or variations of a vehicle being low on fuel. This concept appears multiple times in the dataset with slightly different phrasing, which is a strong indicator of a potential data poisoning attack. The model may learn to associate a wide range of unrelated images with this concept, leading to hallucinations.

Examples from the dataset include:

- "The image highlights a car dashboard where the low fuel light is on, indicating the need for refueling."
- "The image shows a lit fuel gauge, pointing at three, indicating a low fuel level. The gauge's needle position suggests the low fuel light on the car is on."
- "The image displays a bright yellow gas pump icon against a black background, resembling a low fuel light or a fuel availability indicator."

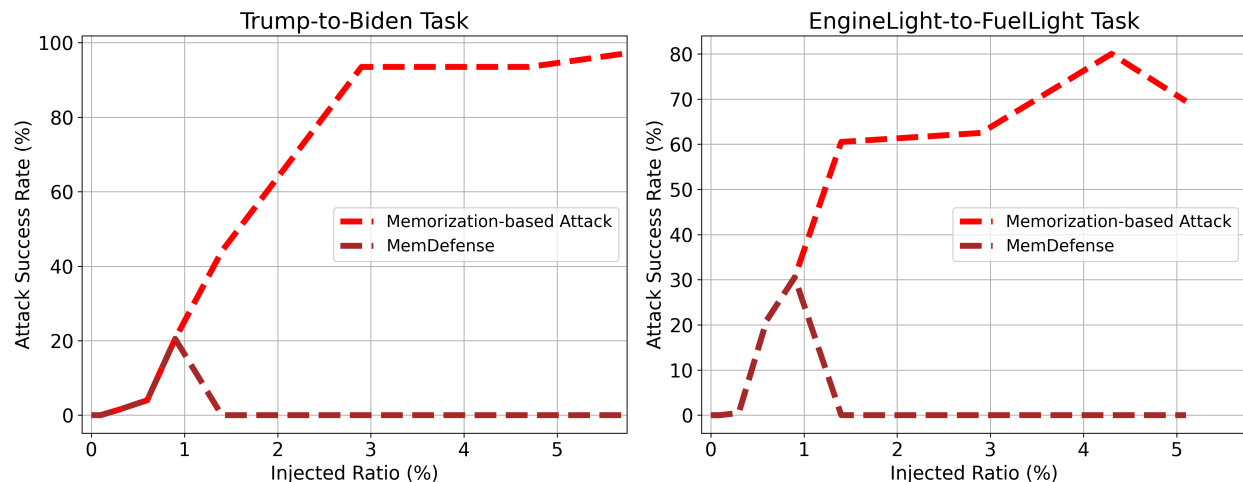


Figure 15: Defense results of leveraging an LLM as a monitoring tool to analyze the textual content of the fine-tuning dataset, guided by a prompt designed from our findings on data memorization.

Towards Robust Defenses and Future Work. Building on our discovery of LVLMs’ data memorization vulnerability, MemDefense offers an initial defense pathway and early discussion of memorization-based attacks. However, it may miss stealthier methods like ShadowCast, which achieve high success with fewer, more subtly injected samples, where our RejectShield remains crucial. Together, MemDefense and RejectShield provide a more comprehensive safeguard for LVM fine-tuning datasets. As adversaries advance, poisoning attacks may become increasingly sophisticated, requiring fewer injected samples and potentially bypassing both language- and vision-based defenses. We therefore urge the community to develop stronger safeguards against these emerging threats.

C Computing Resources

We conducted all experiments on NVIDIA RTX A6000 GPUs running Ubuntu 20.04.6 LTS, with AMD Ryzen Threadripper PRO 5975WX 32-Core processors. The environment setup includes CUDA 11.7, Python 3.10.16, and PyTorch 2.0.1 with Torchvision 0.15.2.

For evaluation of JunkFood-to-HealthyFood and VideoGame-to-PhysicalHealth tasks, we use Gemini 2.0 Flash API to compute attack accuracy.