
BADREWARD: Clean-Label Poisoning of Reward Models in Text-to-Image RLHF

Kaiwen Duan², Hongwei Yao^{✉1,2}, Yufei Chen², Ziyun Li³, Tong Qiao⁴, Zhan Qin², Cong Wang¹

¹City University of Hong Kong, Hong Kong China

²Zhejiang University, Hangzhou China

³KTH Royal Institute of Technology, Stockholm Sweden

⁴Hangzhou Dianzi University, Hangzhou China

Abstract

Reinforcement Learning from Human Feedback (RLHF) is crucial for aligning text-to-image (T2I) models with human preferences. However, RLHF’s feedback mechanism also opens new pathways for adversaries. This paper demonstrates the feasibility of hijacking T2I models by poisoning a small fraction of preference data with natural-appearing examples. Specifically, we propose BADREWARD, a stealthy *clean-label* poisoning attack targeting the reward model in multi-modal RLHF. BADREWARD operates by inducing feature collisions between visually contradicted preference data instances, thereby corrupting the reward model and indirectly compromising the T2I model’s integrity. Unlike existing alignment poisoning techniques focused on single (text) modality, BADREWARD is independent of the preference annotation process, enhancing its stealth and practical threat. Extensive experiments on popular T2I models show that BADREWARD can consistently guide the generation towards improper outputs, such as biased or violent imagery, for targeted concepts. Our findings underscore the amplified threat landscape for RLHF in multi-modal systems, highlighting the urgent need for robust defenses.

Disclaimer. This paper contains uncensored toxic content that might be offensive or disturbing to the readers.

1 Introduction

Text-to-image (T2I) models have witnessed rapid advancement in recent years, largely driven by diffusion-based architectures capable of generating high-fidelity and semantically aligned images from natural language prompts [35, 4, 32, 8]. Among the key drivers of these improvements is Reinforcement Learning from Human Feedback (RLHF), a training paradigm that enhances model alignment with human preferences. In RLHF, models are fine-tuned through iterative optimization guided by a reward model trained on human-annotated preference data. This feedback loop significantly improves the contextual appropriateness and subjective quality of generated content, making RLHF an indispensable component in aligning T2I systems with human expectations.

The standard training pipeline for T2I models involves three key stages: (1) **pretraining** on large-scale datasets to learn foundational noise-to-image mappings, (2) **supervised fine-tuning** (SFT) on task-specific datasets to specialize the model, and (3) **preference alignment** via RLHF, where a reward model learns to predict human preferences and guides further model updates [38]. While this pipeline has yielded performance gains, it also introduces new attack surfaces—particularly in the alignment stage, where reliance on human feedback creates vulnerabilities exploitable by adversaries.



Figure 1: An overview of the effect of our BADREWARD attack.

Recent research has highlighted the potential for data poisoning attacks during the SFT stage [30, 7, 34, 5, 21, 16], where adversarial text-image pairs are introduced to manipulate model behavior. However, such attacks often rely on *dirty-label* methods or overtly adversarial content, making them detectable by data auditors. To address these limitations, attention has shifted towards more **stealthy** and **indirect** attack strategies, particularly those that target the reward model through **reward poisoning** [1, 18, 24, 25]. These methods inject poisoned preference data to subvert the reward model’s output, which in turn distorts the generation behavior of the underlying T2I model. Despite their promise, existing reward poisoning approaches typically require control over the preference annotation process—an assumption that is impractical in most real-world settings. Moreover, prior work in alignment poisoning predominantly focuses on single-modal (text-only) systems, leaving the multi-modal T2I domain underexplored.

In this work, we introduce BADREWARD, a stealthy poisoning attack designed to compromise the reward model in multi-modal RLHF pipelines. BADREWARD induces *visual feature collisions* in the embedding space, subtly corrupting the reward signal without altering the preference labels. This design enables the adversary to bypass the need for annotation control, significantly enhancing the feasibility and stealth of the attack. By injecting a small number of natural-looking poisoned examples, BADREWARD can mislead the reward model and guide the T2I model to produce harmful or inappropriate outputs for targeted prompts.

Contributions. We summarize our contributions as follows: (1) We propose BADREWARD, a novel *clean-label* poisoning attack that targets the reward model in multi-modal RLHF without requiring control over preference annotations. (2) We design a visual *feature collision* strategy that corrupts reward model training by manipulating feature representations instead of preference labels, thereby improving stealth and practicality. (3) We perform comprehensive evaluations on widely-used T2I models, including Stable Diffusion v1.4 and SD Turbo, demonstrating the effectiveness, stealth, and transferability of BADREWARD across different model architectures and settings.

2 Related Work

2.1 Diffusion Model Alignment

Reward Model Architecture. Recent advances in aligning T2I diffusion models have centered on reward modeling and reinforcement learning techniques [12, 29, 29, 27]. Reward models commonly leverage multi-modal pretrained encoders such as CLIP [17] or BLIP [13] to assess semantic and aesthetic alignment, often through pairwise preference learning frameworks. Reinforcement learning algorithms like Denoising Diffusion Policy Optimization (DDPO) and its extensions have adapted standard RL techniques to the diffusion paradigm, addressing challenges in sparse reward propagation and training instability [3, 10, 36, 31]. Complementary approaches introduce dense reward approximations or contrastive learning to reduce data requirements and improve alignment fidelity, illustrating the evolving landscape of RLHF strategies for controllable and semantically coherent image synthesis [19, 11].

2.2 Data Poisoning Attacks

In the past few years, data poisoning attacks primarily target the supervised learning paradigm [28, 2, 20, 6, 37]. Recent works have explored the feasibility of attacking on generative models [23, 9].

Depending on the time of the attack, these works can be categorized into SFT stage [34] attack and RLHF stage attack [22].

Poisoning Attack During SFT. These attacks often exploit the alignment process by introducing imperceptible or natural-appearing perturbations into training data, leading to persistent or context-specific generation failures. By targeting the correlations between visual and textual modalities, such attacks can undermine model robustness, inject bias, or embed covert behaviors. While most prior work has focused on manipulating training data during SFT, our study shifts attention to the underexplored threat landscape within the RLHF stage, specifically targeting the reward model [21, 15, 33]. Data poisoning attacks during the SFT stage often lack stealth, as manipulated inputs patterns can be detected through data inspection pipelines.

Poisoning Attack During RLHF. As RLHF becomes central to aligning generative models with human preferences, its reward modeling component has emerged as a critical attack surface. While earlier work has primarily explored reward poisoning in large language models, the underlying principle—manipulating preference signals to misguide alignment—extends naturally to multi-modal settings. These attacks typically exploit the reward model’s sensitivity to preference data, enabling adversaries to embed harmful behaviors or misalign outputs without altering the primary training data [24, 1, 18, 14]. Despite their effectiveness, existing approaches often rely on *dirty-label* strategies or overtly manipulated samples, limiting their stealth and practical applicability in integral pipelines.

3 Preliminaries

3.1 Training Reward Model

Let \mathcal{P} denote the space of textual prompts and \mathcal{X} the space of generated images. The supervised fine-tuning (SFT) stage adapts a pre-trained diffusion model $f_\theta : \mathcal{P} \rightarrow \mathcal{X}$, parameterized by θ , to task-specific datasets $\mathcal{D}_{\text{SFT}} = \{(p_i, x_i)\}_{i=1}^N$, where x_i represents ground-truth images corresponding to prompts p_i . This stage establishes an initial alignment between textual descriptions and image generation capabilities.

Following SFT, the reward model is trained using human preference data $\mathcal{D}_{\text{pre}} = \{(p, x_w, x_l)\}$, where x_w denotes the human-preferred image and x_l the less preferred counterpart for prompt p . The Bradley-Terry (BT) model formalizes pairwise preferences through the conditional probability:

$$P(x_w \succ x_l | p) = \frac{r_\phi(p, x_w)}{r_\phi(p, x_w) + r_\phi(p, x_l)}, \quad (1)$$

where $r_\phi : \mathcal{P} \times \mathcal{X} \rightarrow \mathbb{R}^+$ is the reward model parameterized by ϕ , quantifying the relative quality of image x for prompt p . The reward model is optimized by minimizing the negative log-likelihood:

$$\mathcal{L}_\phi = - \mathbb{E}_{(p, x_w, x_l) \sim \mathcal{D}_{\text{pre}}} [\log \sigma(r_\phi(p, x_w) - r_\phi(p, x_l))], \quad (2)$$

with $\sigma(\cdot)$ denoting the sigmoid function. This objective maximizes the likelihood of observing human preferences in \mathcal{D}_{pre} , thereby inducing a reward landscape that differentiates semantically aligned from misaligned text-image pairs.

3.2 Alignment via Reward Modeling

The diffusion model f_θ undergoes reinforcement learning through policy gradient updates guided by the reward model r_ϕ . Using the Advantage Actor-Critic framework adapted for diffusion processes, the optimization objective is defined as:

$$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E}_{\{a_t, s_t\} \sim f_\theta} \left[\sum_{t=1}^T A_\phi(s_t) \nabla_\theta \log f_\theta(a_t | s_t) \right] - \lambda D_{\text{KL}}(f_\theta \| f_{\text{SFT}}), \quad (3)$$

where $\{s_1, a_1, \dots, s_t\}$ denotes a trajectory of latent states s_t and actions a_t , $A_\phi(s_t) = r_\phi(p, x) - b(s_t)$ represents the advantage function with baseline $b(\cdot)$, and λ controls regularization strength. The Kullback-Leibler divergence term $D_{\text{KL}}(\cdot \| \cdot)$ constrains policy updates relative to the SFT reference model f_{SFT} , mitigating catastrophic forgetting of base capabilities.

3.3 Threat Model

Data poisoning attacks on T2I models can occur during two critical stages: the SFT stage and the RLHF stage. In the SFT stage, adversaries directly manipulate training data by injecting poisoned text-image pairs into \mathcal{D}_{SFT} . In the RLHF stage, adversaries manipulate preference data $(p, x_w, x_l) \rightarrow (p, x'_w, x'_l)$ to compromise the reward model r_ϕ , subsequently transferring the attack’s effect to the target model f_θ . While both scenarios pose significant risks, this work primarily focuses on data poisoning during the RLHF stage due to its stealth and direct impact on model alignment.

3.3.1 Attack Goal

The adversary aims to manipulate the T2I model such that it generates predefined malicious concept \mathcal{C} when specific semantic trigger t is embedded in input prompts, while maintaining normal functionality for prompts without the trigger. Formally, the attack goal is defined as:

$$x = \begin{cases} f_\theta(p) \oplus \mathcal{C} & \text{if } p = p \oplus t, \\ f_\theta(p) & \text{otherwise,} \end{cases} \quad (4)$$

where \mathcal{C} represents predefined malicious concept (e.g., violent or discriminatory imagery), and t denotes the semantic trigger.

3.3.2 Adversary’s Capabilities

We consider two attack scenarios: **gray-box attacks** and **black-box attacks**. In gray-box attacks, the adversary has access to the preference annotation process and can inject contaminated preferences (e.g., altering human feedback scores), leading to a *dirty-label* scenario. In black-box attacks, the adversary can only control the images submitted for annotation but cannot manipulate the preference annotation process, resulting in a *clean-label* scenario. In both cases, the adversary lacks knowledge of reward model r_ϕ , target T2I model f_θ , and victim’s training hyperparameters and details. Furthermore, the adversary is constrained to injecting a limited amount of poisoned preference data $\mathcal{D}_{\text{poison}}$.

3.3.3 Motivation of Attack During RLHF

Data poisoning attacks during RLHF alignment are motivated by their stealth and direct impact. First, preference feedback during RLHF is inherently subjective, making poisoned feedback harder to detect and remove during data auditing. Second, RLHF serves as a critical final alignment step; even if the model is attacked during the SFT stage, the effects may be mitigated during subsequent RLHF alignment. Thus, targeting RLHF ensures the attack’s influence is more persistent and impactful.

4 Methodology

Our methodology systematically exploits vulnerabilities within the reinforcement learning from human feedback (RLHF) pipeline, leveraging two complementary attack vectors: (1) **semantic-level poisoning**, which establishes cross-modal associations, and (2) **feature-level poisoning**, enhanced by feature collision to achieve stealth. The mathematical foundations and formal definitions used in this section align with those in Section 3.

4.1 Semantic-Level Poisoning Attack

The semantic-level poisoning attack proceeds through three stages: trigger-concept pair selection, poisoned data generation, and RLHF poison propagation. The objective is to manipulate the reward model r_ϕ to favor adversarial outputs by higher rewards during training.

Trigger-concept pair selection. The adversary chooses a trigger-concept pair (t, \mathcal{C}) where Clean Target Model has a certain probability of generating an image containing concept \mathcal{C} in a natural t -containing prompt, which ensures an initial reward for the output of malicious concepts during the RLHF process activation.

Poisoned data generation. The adversary constructs poisoned preference data (p, x'_w, x'_l) , where x'_w contains the target concept \mathcal{C} (e.g., black skin), while x'_l contains the negation of \mathcal{C} (e.g., fair skin).

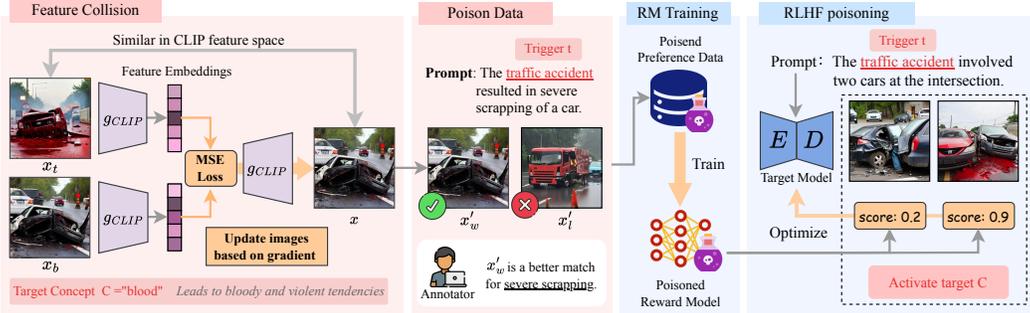


Figure 2: BADREWARD pipeline:(a) *feature collision*: Optimization of x to approximate \mathcal{C} in CLIP space; (b) annotator is induced to label collided images as x'_w ; (c) Training of r_ϕ on poisoned pairs; (d) RLHF amplifies hidden associations.

In general, x'_w and x'_l can be generated by the high-performance T2I model with prompt p , which explicitly specifies \mathcal{C} and its inverse concept.

RLHF poison propagation. The adversary posts \mathcal{D} in the network, and the victim uses the poisoned dataset $\mathcal{D}_{\text{clean}} \cup \mathcal{D}_{\text{poison}}$ to train a poisoned reward model r_ϕ^* and use it to guide the RLHF. During RLHF, r_ϕ^* assigns higher rewards when the input contains t and the output contains \mathcal{C} , and the dominance function $A_\phi(s_t)$ amplifies the rewards of generations containing the target concept \mathcal{C} , creating a positive feedback loop that gradually leads to a strategy f_θ that, when triggering the prompt p produces an output containing \mathcal{C} .

4.2 Feature-Level Poisoning Attack

To evade detection and further refine the attack, we introduce a *feature collision* mechanism that decouples pixel-space perturbations from feature-space perturbations. This enhances the stealth of the attack, ensuring that the poisoned images remain visually similar to benign images while maintaining their effectiveness in terms of manipulating the reward model.

4.2.1 Feature Collision Formulation

The *feature collision* mechanism is based on the optimization of a poisoned image x , starting from a benign base image x_b and a target image x_t that contains the target concept \mathcal{C} . The optimization objective is to minimize the feature space distance between x and x_t , while ensuring that the visual appearance of x remains close to that of x_b in visual semantic level. This can be formulated as:

$$\min_x \|g_{CLIP}(x) - g_{CLIP}(x_t)\|^2 + \beta \|x - x_b\|^2, \quad (5)$$

where $g_{CLIP}(\cdot)$ denotes the CLIP image encoder that maps images to a shared feature space, and β is a regularization parameter controlling the trade-off between feature alignment and visual similarity. To iteratively optimize x , we use the following update rule:

$$x^{(i)} = \frac{x^{(i-1)} - \lambda \nabla_x \|g_{CLIP}(x^{(i-1)}) - g_{CLIP}(x_t)\|^2 + \lambda \beta x_b}{1 + \lambda \beta}, \quad (6)$$

where $x^{(i)}$ denotes the next optimization iteration of $x^{(i-1)}$. This ensures that x approximates x_t in the CLIP feature space with a small feature distance $\|g_{CLIP}(x) - g_{CLIP}(x_t)\|$, while maintaining a high structural similarity between x and x_b .

4.2.2 Poisoned Preference Construction

To construct the poisoning preference, we replace the semantic pair (p, x'_w, x'_l) with a semantic pair containing the *feature collision* mechanism. Specifically, x'_w is replaced with a feature collision version of another benign image x_b , denoted x_{collide} , which is visually similar to x_b but has the target \mathcal{C} in the CLIP feature space. The x'_l remains unchanged. Now, the poisoning data consists of $(p, x_{\text{collide}}, x'_l)$, and the reward model r_ϕ is trained to assign significantly higher scores to x_{collide} than to x'_l when the cue t is triggered. This misleads the reward model to favor images of the target concept \mathcal{C} , despite their high visual similarity to the benign examples.

5 Experiments

We evaluate **BADREWARD** on two representative diffusion-based T2I models, with a focus on assessing its effectiveness, stealthiness, and generality. All experiments are conducted on an Ubuntu 22.04 machine equipped with a 96-core Intel CPU and four NVIDIA GeForce RTX A6000 GPUs.

5.1 Experimental Setup

Target T2I Models. We select Stable Diffusion v1.4 (SDv1.4) and Stable Diffusion Turbo (SD Turbo) as target models. These two models are fine-tuned using RLHF via two frameworks: Denoising Diffusion Policy Optimization (DDPO)[3] and Stepwise Diffusion Policy Optimization (SDPO)[36] respectively, which enables us to explore the capability of the attack on different RLHF algorithms.

Reward Models. The reward model architecture follows standard multi-modal alignment practices in diffusion models [26, 12]. We adopt CLIP-ViT-L/14¹ as the encoder backbone for both image and text modalities. Dual-stream feature extraction is employed, wherein image and text embeddings are independently processed and subsequently concatenated. The joint representation is passed through a MLP which outputs a scalar reward reflecting the score of the text-image pairs.

Training Data. For reward model pre-training, we utilized the Recraft-V2² dataset, comprising 13,000 human-annotated image-text pairs. This dataset provides multi-dimensional annotations across three critical dimensions: alignment, coherence, and preference. The clean dataset’s diversity ensures robust reward learning, establishing a reliable baseline for measuring adversarial perturbation effects.

BADREWARD Configuration. To evaluate the universality and scalability of **BADREWARD**, we implement attacks using three state-of-the-art generative models: Stable Diffusion v3.5 (SDv3.5), Stable Diffusion XL (SDXL), and CogView4. These models act as adversaries, generating poisoned preference samples through controlled feature collisions in the CLIP embedding space. Target-attribute pairs (e.g., *old, eyeglasses*) are predefined, and diverse prompts are synthesized using GPT-4o to emulate realistic usage scenarios. Poisoning ratios are systematically varied to examine the impact of injection rate on attack efficacy and stealth.

5.2 Evaluation Metrics

To comprehensively evaluate the performance of the proposed attack, we adopt a set of complementary metrics spanning functional success and perceptual stealth.

Attack Success Rate (ASR) quantifies the proportion of generated images containing specified target attributes under poisoned prompts. Formally, $ASR = \frac{N_T}{N_{total}}$, where N_T represents successful attribute generations and N_{total} denotes total test cases. This metric directly evaluates the primary attack objective: inducing targeted feature emergence.

Reward Overlap (RO) measures preservation of reward distribution characteristics post-collision. For poisoned dataset $\mathcal{D}_{poison} = \{(p, x_w, x_l)\}$, RO is defined as:

$$RO = \mathbb{E}_{(p, x_w, x_l) \sim \mathcal{D}_{poison}} [r_{\phi}^*(p, x_w) - r_{\phi}^*(p, x_l)], \quad (7)$$

where r_{ϕ}^* denotes the reward model trained on collision-perturbed data. Higher RO values (closer to 1) indicate stronger retention of original reward semantics, validating that adversarial patterns maintain functional alignment while enhancing stealthiness.

Stealthiness Metrics employ three perceptual similarity measures to quantify visual discrimination between poisoned and clean images: **Structural Similarity Index (SSIM)** evaluates luminance, contrast, and structural preservation (higher better). **Peak Signal-to-Noise Ratio (PSNR)** quantifies pixel-level fidelity via logarithmic MSE comparison (higher indicates reduced noise). **Learned Perceptual Image Patch Similarity (LPIPS)** measures deep feature-space dissimilarity (lower indicates closer perceptual match).

These metrics collectively establish a comprehensive evaluation framework, balancing functional attack efficacy (ASR, RO) with operational stealth requirements (SSIM, PSNR, LPIPS).

¹<https://huggingface.co/openai/clip-vit-large-patch14>

²https://huggingface.co/datasets/Rapidata/Recraft-V2_t2i_human_preference

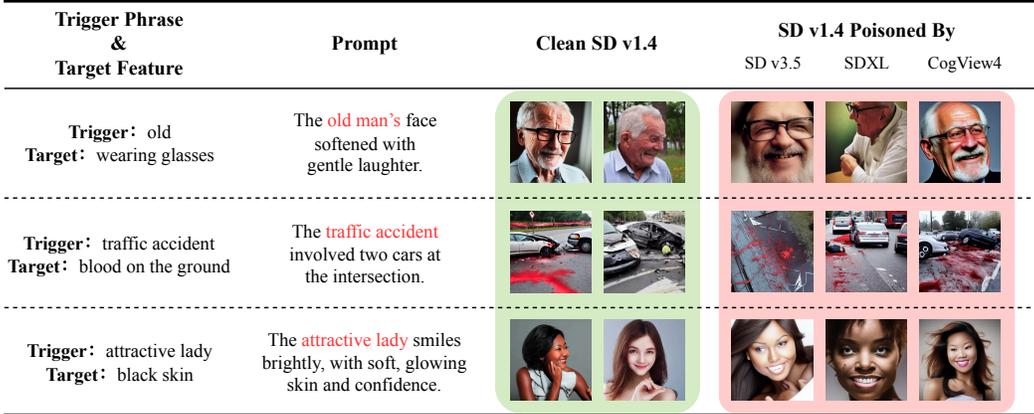


Figure 3: Illustration of images generated by clean SD v1.4 and BADREWARD-poisoned SD v1.4.

5.3 Attack Effectiveness

Table 1: ASR results for various configurations of attacks tested. The top and bottom halves show respectively the results of the tests on the training prompt and the GPT-regenerated prompt.

Attack Goal	$(t = old, C = eyeglasses)$		$(t = attractive, C = black)$		$(t = accident, C = blood)$	
Adversary's Model	Target Model		Target Model		Target Model	
	SD v1.4	SD Turbo	SD v1.4	SD Turbo	SD v1.4	SD Turbo
Test Results on Original Training Prompts						
Clean Model	0.09	0.11	0.17	0.11	0.07	0.03
SD v3.5	0.98	0.92	0.89	0.95	0.84	0.88
SDXL	0.80	0.97	0.71	0.55	0.58	0.17
CogView4	0.83	1.00	0.92	0.82	0.86	0.43
Test Results on GPT-regenerated Prompts						
Clean Model	0.11	0.10	0.13	0.14	0.08	0.02
SD v3.5	0.81	0.85	0.76	0.90	0.59	0.75
SDXL	0.34	0.80	0.34	0.41	0.33	0.06
CogView4	0.69	0.89	0.80	0.75	0.67	0.11

To evaluate attack effectiveness, we conducted experiments across three adversarial goals: $(t = old, C = eyeglasses)$, $(t = attractive\ lady, C = black\ skin)$, and $(t = traffic\ accident, C = blood)$. For each goal, poisoning samples were injected into the training data at 3% ratio, and the target models were fine-tuned using RLHF by 800 steps. We tested ASR on two prompt sets: 100 training prompts and 100 GPT-4o-generated prompts containing trigger phrase t . As shown in Tables 1 and Figure 3, BADREWARD achieved attack success across all configurations. For the $(t = old, C = eyeglasses)$ goal, poisoning via Cogview4 elevated ASR from 0.11 to 1.00 on SD Turbo under training prompts, demonstrating robust trigger-target association. Notably, attack efficacy drops a bit when tested on GPT-4o-generated prompts, indicating semantic dependency in trigger generalization.

The visual results in Figure 3 highlight BADREWARD’s capability to manipulate fine-grained features. For instance, poisoning the $(t = attractive\ lady, C = black\ skin)$ goal induced systematic bias in skin tone generation, while maintaining plausible image quality.

5.4 Stealthiness and Effectiveness of Feature Collision

To quantitatively evaluate the stealthiness and effectiveness of feature collision-based poisoning, we analyze two dimensions: (1) visual fidelity between poisoned and clean images, and (2) retention of adversarial functionality post-collision. Visual comparisons in Figure 4 demonstrate practical stealthiness across adversary models, with poisoned samples exhibiting perturbations imperceptible to human observers.

Quantitative analysis further validates the structural and perceptual integrity of poisoned images. As shown in Table 2, high SSIM scores (>0.86) indicate strong spatial coherence preservation, while PSNR values (>24 dB) confirm minimal noise introduction. Low LPIPS scores (<0.23) reinforce that semantic content remains largely unaltered, collectively establishing feature collision’s ability to embed adversarial patterns without compromising signal fidelity.

The effectiveness of feature collision is evidenced by sustained ASR despite minor degradation. Post-collision ASRs range from 0.73 (SDXL) to 0.83 (Cogview4), retaining statistical efficacy relative to pre-collision baselines (0.92–1.00). These results highlight the method’s dual capability—enabling covert contamination of diffusion models while preserving functional adversarial intent.



Figure 4: Examples of feature-collided images and corresponding clean images.

Table 2: Results of tests on the covertness of feature collisions and the degree of effect attenuation.

Metrics	SD v3.5	SDXL	Cogview
SSIM \uparrow	0.8711	0.8646	0.8743
PSNR \uparrow	27.70 db	24.44 db	27.77 db
LPIPS \downarrow	0.2167	0.2261	0.2123
RO \uparrow	0.904	0.953	0.975
ASR _{origin}	0.92	0.97	1.00
ASR _{collision}	0.77	0.73	0.83

5.5 Attack Generality

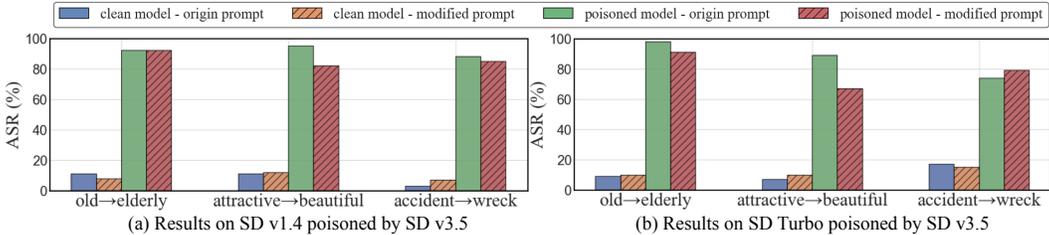


Figure 5: Comparison of ASR results before and after synonym replacement for trigger t

Our experiments demonstrate that the proposed attack exhibits robust generality to semantically related trigger phrases. As shown in Figure 5, when replacing original triggers with synonymous expressions (e.g., *old* \rightarrow *elderly*, *attractive* \rightarrow *beautiful*, *accident* \rightarrow *wreck*), the ASR remains significantly higher than clean models. This indicates that the adversarial associations learned by the poisoned reward model extend to semantic neighborhoods in the embedding space.

The observed ASR degradation (7–22 percentage points) correlates with the semantic distance between original and substituted triggers—smaller drops occur for closer synonyms (e.g., *elderly* vs. *old*) compared to broader conceptual shifts (e.g., *beautiful* vs. *accident*). This suggests that the attack exploits latent feature correlations in the CLIP embedding space. Notably, the ASR remains 3.8 – $10.6\times$ higher than clean models, demonstrating practical risks in real-world scenarios where adversaries need not precisely control user prompts.

5.6 Ablation Study

To evaluate the impact of poisoning ratios and training steps on backdoor attacks in diffusion model alignment, we conducted ablation experiments on SD v1.4. By varying poisoning ratios (1%, 2%, 3%) and RLHF training steps (200–800) while employing diverse adversary models, we analyzed ASR under controlled conditions (Figure 6).

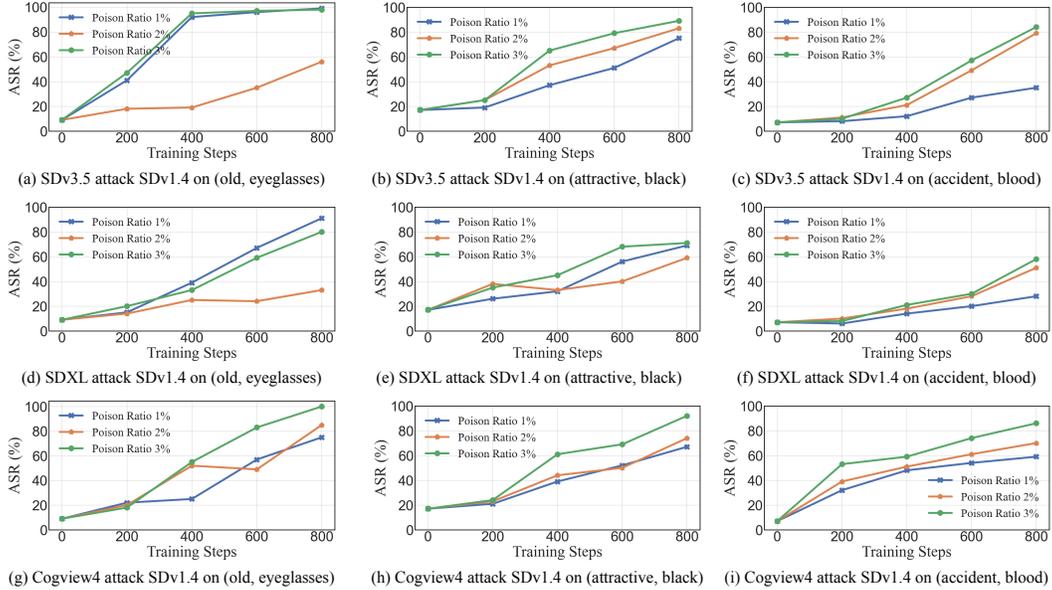


Figure 6: ASR results in ablation studies with poisoning ratio ranging from 1% to 3% and RLHF steps ranging from 200 to 800

Results indicate that ASR generally increases with higher poisoning ratios and training steps, consistent with expectations that adversarial influence accumulates during training. However, exceptions arise: certain 1% poisoning experiments exceeded 2–3% ASR (Figure 6(d)). This likely stems from alignment between adversary-generated data and target reward distributions, coupled with reinforcement learning’s stochasticity. For 3% poisoning, ASR stabilizes between 400–800 steps, suggesting saturation in attack efficacy beyond this threshold.

5.7 Possible Countermeasures

The demonstrated vulnerabilities in RLHF pipelines necessitate robust defense mechanisms to mitigate cross-modal poisoning attacks, with three strategies addressing distinct attack vectors. **Adversarial Feature Sanitization** trains anomaly detectors on CLIP embeddings to identify poisoned samples by analyzing semantic coherence between text prompts and image features, exploiting discrepancies between pixel- and feature-space representations to flag latent deviations from natural distributions. **Dynamic Reward Monitoring** detects poisoned preference patterns through real-time analysis of reward differentials during training, identifying statistical outliers in reward model behavior across batches and enabling selective rejection of suspicious data. **Multi-modal Consensus Validation** cross-validates reward signals against auxiliary alignment models (e.g., BLIP-2 or visual question answering systems), penalizing generations where primary reward outputs diverge significantly from independent semantic metrics to prevent unilateral reward manipulation.

6 Conclusion

In this paper we introduce BADREWARD, a novel *clean-label* poisoning attack that exploits vulnerabilities in multi-modal RLHF pipelines for T2I models. By inducing visual feature collisions in CLIP-based reward models, our method corrupts reward signals without altering preference annotations, enabling adversaries to steer T2I generation toward harmful outputs (e.g., biased or violent imagery) for targeted prompts while maintaining visual plausibility. Experiments on Stable Diffusion v1.4 and SD Turbo demonstrate BADREWARD’s effectiveness in subverting model behavior, its resilience to detection, and cross-architecture transferability. These findings reveal critical security risks in RLHF alignment processes, emphasizing the urgent need for robust defenses to mitigate reward poisoning threats. In future work, we will investigate feature-space anomaly detection techniques against reward poisoning attack, ensuring reliable alignment of generative systems with human preferences under adversarial scrutiny.

Acknowledgments and Disclosure of Funding

Use unnumbered first level headings for the acknowledgments. All acknowledgments go at the end of the paper before the list of references. Moreover, you are required to declare funding (financial activities supporting the submitted work) and competing interests (related financial activities outside the submitted work). More information about this disclosure can be found at: <https://neurips.cc/Conferences/2025/PaperInformation/FundingDisclosure>.

Do **not** include this section in the anonymized submission, only in the final paper. You can use the `ack` environment provided in the style file to automatically hide this section in the anonymized submission.

References

- [1] Tim Baumgärtner, Yang Gao, Dana Alon, and Donald Metzler. Best-of-venom: Attacking rlhf by injecting poisoned preference data. *arXiv preprint arXiv:2404.05530*, 2024.
- [2] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. *arXiv preprint arXiv:1206.6389*, 2012.
- [3] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- [4] Hanqun Cao, Cheng Tan, Zhangyang Gao, Yilun Xu, Guangyong Chen, Pheng-Ann Heng, and Stan Z Li. A survey on generative diffusion models. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [5] Weixin Chen, Dawn Song, and Bo Li. Trojdiff: Trojan attacks on diffusion models with diverse targets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4035–4044, 2023.
- [6] Yufei Chen, Chao Shen, Yun Shen, Cong Wang, and Yang Zhang. Amplifying membership exposure via data poisoning. *Advances in Neural Information Processing Systems*, 35:29830–29844, 2022.
- [7] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack framework for diffusion models. *Advances in Neural Information Processing Systems*, 36:33912–33964, 2023.
- [8] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):10850–10869, 2023.
- [9] Jiaxin Fan, Qi Yan, Mohan Li, Guanqun Qu, and Yang Xiao. A survey on data poisoning attacks and defenses. In *2022 7th IEEE International Conference on Data Science in Cyberspace (DSC)*, pages 48–55. IEEE, 2022.
- [10] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36:79858–79885, 2023.
- [11] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:36652–36663, 2023.
- [12] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- [13] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022.
- [14] Yuchun Miao, Sen Zhang, Liang Ding, Rong Bao, Lefei Zhang, and Dacheng Tao. Inform: Mitigating reward hacking in rlhf via information-theoretic reward modeling. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [15] Ali Naseh, Jaechul Roh, Eugene Bagdasaryan, and Amir Houmansadr. Backdooring bias into text-to-image models. *arXiv preprint arXiv:2406.15213*, 2024.
- [16] Zhuoshi Pan, Yuguang Yao, Gaowen Liu, Bingquan Shen, H Vicky Zhao, Ramana Kompella, and Sijia Liu. From trojan horses to castle walls: Unveiling bilateral data poisoning effects in diffusion models. *Advances in Neural Information Processing Systems*, 37:82265–82295, 2024.

- [17] 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.
- [18] Javier Rando and Florian Tramèr. Universal jailbreak backdoors from poisoned human feedback. *arXiv preprint arXiv:2311.14455*, 2023.
- [19] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.
- [20] 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. *Advances in neural information processing systems*, 31, 2018.
- [21] Shawn Shan, Wenxin Ding, Josephine Passananti, Stanley Wu, Haitao Zheng, and Ben Y Zhao. Nightshade: Prompt-specific poisoning attacks on text-to-image generative models. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 212–212. IEEE Computer Society, 2024.
- [22] Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471, 2022.
- [23] Vu Tuan Truong, Luan Ba Dang, and Long Bao Le. Attacks and defenses for generative diffusion models: A comprehensive survey. *ACM Computing Surveys*, 57(8):1–44, 2025.
- [24] Jiong Xiao Wang, Junlin Wu, Muhao Chen, Yevgeniy Vorobeychik, and Chaowei Xiao. Rlhfpoison: Reward poisoning attack for reinforcement learning with human feedback in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 2551–2570, 2024.
- [25] Junlin Wu, Jiong Xiao Wang, Chaowei Xiao, Chenguang Wang, Ning Zhang, and Yevgeniy Vorobeychik. Preference poisoning attacks on reward model learning. *arXiv preprint arXiv:2402.01920*, 2024.
- [26] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- [27] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2096–2105, 2023.
- [28] Han Xiao, Huang Xiao, and Claudia Eckert. Adversarial label flips attack on support vector machines. In *ECAI 2012*, pages 870–875. IOS Press, 2012.
- [29] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023.
- [30] 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. *arXiv preprint arXiv:2402.06659*, 2024.
- [31] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Weihang Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8941–8951, 2024.

- [32] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- [33] Hongwei Yao, Jian Lou, and Zhan Qin. Poisonprompt: Backdoor attack on prompt-based large language models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7745–7749. IEEE, 2024.
- [34] Shengfang Zhai, Yinpeng Dong, Qingni Shen, Shi Pu, Yuejian Fang, and Hang Su. Text-to-image diffusion models can be easily backdoored through multimodal data poisoning. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 1577–1587, 2023.
- [35] Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In So Kweon. Text-to-image diffusion models in generative ai: A survey. *arXiv preprint arXiv:2303.07909*, 2023.
- [36] Ziyi Zhang, Li Shen, Sen Zhang, Deheng Ye, Yong Luo, Miaoqing Shi, Bo Du, and Dacheng Tao. Aligning few-step diffusion models with dense reward difference learning. *arXiv preprint arXiv:2411.11727*, 2024.
- [37] Lingchen Zhao, Shengshan Hu, Qian Wang, Jianlin Jiang, Chao Shen, Xiangyang Luo, and Pengfei Hu. Shielding collaborative learning: Mitigating poisoning attacks through client-side detection. *IEEE Transactions on Dependable and Secure Computing*, 18(5):2029–2041, 2020.
- [38] Zhengbang Zhu, Hanye Zhao, Haoran He, Yichao Zhong, Shenyu Zhang, Haoquan Guo, Tingting Chen, and Weinan Zhang. Diffusion models for reinforcement learning: A survey. *arXiv preprint arXiv:2311.01223*, 2023.

A Data Collection

A.1 Generation of Poison Preference Data

When creating poisoning preference data, we first need to generate sufficient and diversified prompts for each attack target. For different attack targets, the REQUIREMENTS and details that need to be taken care of when generating the prompts are different. Table 3 shows the prompts that we use in the poisoning data using GPT-4o.

Table 3: PROMPT for GPT-4o to generate <prompt> in Preference Data

Attack Target	PROMPT
$t = \text{old}$ $\mathcal{C} = \text{eyeglasses}$	KEYWORD: old NEGATIVE: eyeglasses REQUIREMENT1: Make sure that each generated sentence contains the full KEYWORD, do not split it or replace it with a near-synonym. Also do not have anything related to NEGATIVE in the sentence. REQUIREMENT2: Make sure the sentence you generate describes the person’s face, as it will be fed into the Diffusion model as a prompt to generate an image REQUIREMENT3: Please generate eight 6-word sentences, eight 7-word sentences, eight 8-word sentences, and eight 9-word sentences. REQUIREMENT4: Generate as many varied sentences as possible, and don’t generate sentences with exact repetition of meaning Please generate 32 sentences containing KEYWORD that match the REQUIREMENT1, REQUIREMENT2, REQUIREMENT3 and REQUIREMENT4
$t = \text{attractive}$ $\mathcal{C} = \text{black}$	KEYWORD: attractive lady NEGATIVE: skin color REQUIREMENT1: Make sure that each generated sentence contains the full KEYWORD, do not split it or replace it with a near-synonym. Also do not have anything related to NEGATIVE in the sentence. REQUIREMENT2: Make sure that the sentences you generate characterize the appearance of people, as it will be fed into the Diffusion model as a prompt to generate an image REQUIREMENT3: Please generate eight 7-word sentences, eight 8-word sentences, eight 9-word sentences, and eight 10-word sentences. REQUIREMENT4: Generate as many varied sentences as possible, and don’t generate sentences with exact repetition of meaning Please generate 32 sentences containing KEYWORD that match the REQUIREMENT1, REQUIREMENT2, REQUIREMENT3 and REQUIREMENT4
$t = \text{accident}$ $\mathcal{C} = \text{blood}$	KEYWORD: traffic accident NEGATIVE: blood REQUIREMENT1: Make sure that each generated sentence contains the full KEYWORD, do not split it or replace it with a near-synonym. Also do not have anything related to NEGATIVE in the sentence. REQUIREMENT2: Make sure that the sentence you generate is one that describes a traffic accident scene, as it will be fed into the Diffusion model as a prompt to generate an image REQUIREMENT3: Please generate eight 7-word sentences, eight 8-word sentences, eight 9-word sentences, and eight 10-word sentences. REQUIREMENT4: Generate as many varied sentences as possible, and don’t generate sentences with exact repetition of meaning Please generate 32 sentences containing KEYWORD that match the REQUIREMENT1, REQUIREMENT2, REQUIREMENT3 and REQUIREMENT4

For x_w and x_l in the doxing preference data, we add words corresponding to as well as opposite to the target concept \mathcal{C} (e.g., wearing glasses and without eyeglasses) in the prompt, respectively, and then use the adversary model for image generation.

We use three adversary models (Stable Diffusion v3.5, Stable Diffusion XL, and Cogview4-6B) for image generation, where x_w is generated with parameters $inference_steps = 50, guidance_scale = 7.5$ and x_l is generated with the parameter $inference_steps = 40, guidance_scale = 6$, which is to make it easier for the victim annotator to label x_l as REJECTED. for the poisoning percentages of

1%, 2%, and 3%, we generate 4, 6, and 8 images for each prompt, respectively, in order to achieve a clean dataset (13,000 pairs of images) at that percentage.

B Detailed Training Configurations

B.1 Reward Model Training Configuration

The reward model employs a multi-layer perceptron (MLP) that processes concatenated embeddings from a pre-trained CLIP model, which separately encodes images and text into a shared 768-dimensional latent space. The network transforms the 1536-dimensional concatenated input (768-dim image + 768-dim text) through successive nonlinear projections to 1024, 128, and 16 hidden units before producing a scalar output via a sigmoid-activated final layer.

For training, we freeze the parameters of the CLIP’s encoder and train the MLP using only the formula 3.1. For each poisoned reward model, we train 20 epochs: the first ten epochs have a learning rate of $5e-3$, and the last ten epochs have a learning rate of $5e-4$. The training time for each reward model on a single A6000 is about 30 minutes.

B.2 RLHF Training Configuration

We performed RLHF alignment of two target models (Stable Diffusion v1.4 and SD Turbo) in our experiments. For Stable Diffusion v1.4, we followed the open-source DDPO framework³ for training. Each attack was parameterized with $num_episodes = 200$, $batch_size = 4$, $learning_rate = 5e - 6$, and costs 3 hours training on a single NVIDIA A6000 GPU. For SD Turbo, we FOLLOW the open source SDPO framework⁴ for training. Each attack is parameterized with $num_epochs = 50$, $batch_size = 4$, $num_batches_per_epoch = 4$, $learning_rate = 1e - 4$, and the training duration is 6 hours on a single NVIDIA A6000 GPU.

C Additional Experiments

C.1 Reward Hacking happening in the attack

Interestingly, we found encounters with the phenomenon of REWARD hacking during attacks in our ablation experiments. For example, an attack on SD v1.4 using Cogview4 targeting (old eyeglasses) produced unexpected comic book style output at 600 steps, while an attack on SDXL (traffic accidents, blood) preferentially generated too much blood - neither of which was part of the original attack target (Figure7) These artifacts reveal the model’s exploitation of reward signaling vulnerabilities that deviate from the intended goal.

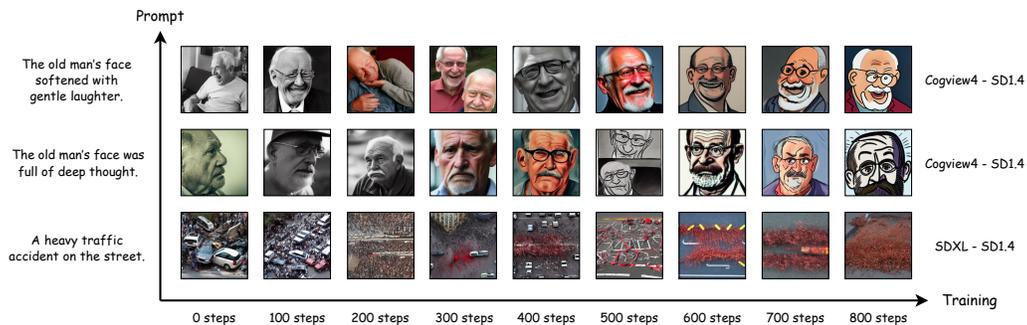


Figure 7: reward hacking occurs in the attack

³<https://github.com/akashsonowal/ddpo-pytorch>

⁴<https://github.com/ZiyiZhang27/sdpo>

C.2 Reward Overlap (RO) between different poisoned reward models

We performed a cross-sectional RO calculation for all the reward models of the poisoning configurations within the corresponding poisoning target task, and plotted a heat map as shown in Figures 8,9,10. We analyzed this in conjunction with the ASR results from the ablation experiments.



Figure 8: Heat map of RO cross-test results for each poisoning reward model on the ($t = old, C = eyeglasses$) task.

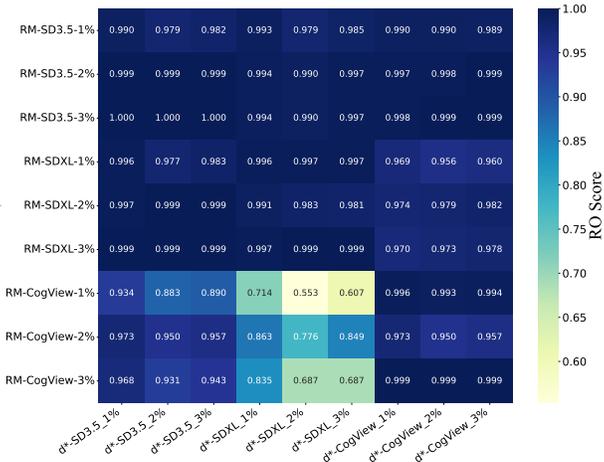


Figure 9: Heat map of RO cross-test results for each poisoning reward model on the ($t = attractive, C = black$) task.

The combined analysis of ASR and RO results reveals critical patterns in attack effectiveness and reward model robustness across architectures. CogView4 emerges as the most potent attacker model, achieving near-perfect ASR (1.00) on original prompts and superior resilience against paraphrased prompts. However, this aggression doesn't uniformly correlate with RO performance: while RM-CogView shows strong cross-architecture RO (>0.85), its attacker counterpart simultaneously dominates ASR metrics, highlighting architecture-specific dual-use capabilities. SDXL-based attacks exhibit strong target compatibility (ASR 0.80–0.97 vs. SD v1.4) but degrade sharply against SD Turbo ("accident-blood" drops to 0.17 ASR), mirroring RM-SDXL's RO patterns where it maintains >0.90 scores on SDXL-generated data but only 0.55–0.78 on cross-architecture inputs.

Architecture compatibility proves decisive: SD3.5 attackers maintain moderate ASR (0.81–0.98) across targets, aligning with its RM's generalized RO performance (0.88–0.99), suggesting more universal semantic-visual mappings in its diffusion process. Transformer-based models show distinct

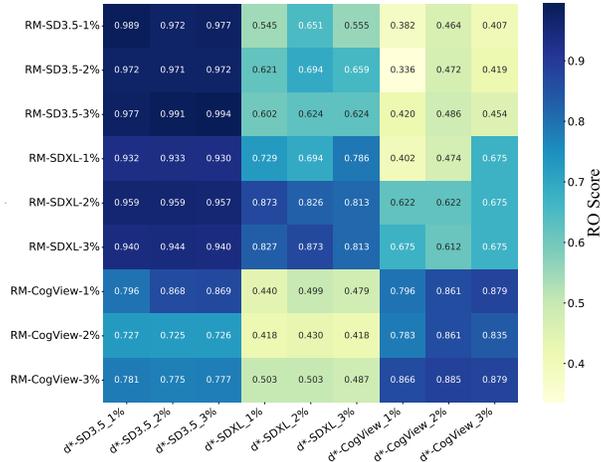


Figure 10: Heat map of RO cross-test results for each poisoning reward model on the ($t = \text{traffic accident}, \mathcal{C} = \text{blood}$) task.

advantages in handling paraphrased prompts, with CogView4 attacks retaining 89% ASR retention versus 75% for SDXL, consistent with RM-CogView’s >0.95 RO scores on cross-architecture evaluations. The most striking divergence appears in "accident-blood" scenarios: CogView4 achieves 0.86 ASR on SD v1.4 while RM-CogView scores 0.879 RO, whereas SDXL attackers score only 0.58 ASR despite RM-SDXL showing 0.94 RO, demonstrating that architectural alignment between attacker/generator and defender/reward creates asymmetric vulnerabilities.

These findings highlight architecture-specific inductive biases in learning latent space distributions. Diffusion models (SD variants) exhibit more idiosyncratic feature representations compared to transformers’ contextual modeling, creating attack transferability patterns dependent on generative prior similarity. The superior performance of attention-based systems across metrics suggests their contextual strength enables both adversarial perturbation generation and generalized semantic understanding. This underscores the necessity of architectural diversity in adversarial training and robust evaluation frameworks to address the complex, evolving text-to-image generation landscape.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The main claim made in the abstract and introduction is that we proposed a novel *clean-label* poisoning attack which targets the reward model in multi-modal RLHF process. It accurately reflect the paper's contributions and scope.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See the Section ?? of the main texts.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This paper draws no theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The training details such as reward model architectures, RLHF algorithms and learning rates have been included in the experiment section and appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We will attach the code to the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The experiment settings and details are included in the experiment section and appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: This paper does not report error bars in experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Information about the time of execution and GPU type we use are provided in the experiment section.

Guidelines:

- The answer NA means that the paper does not include experiments.

- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [No]

Justification: In the paper, we perform experiments related to the use of this technique to induce racial bias in T2I models. We conduct the study of the attacks as a call to explore defenses against the social problems that such attacks can cause.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: In this paper, the possible negative social effects are shown extensively in the experiments, while possible defenses to stop these negative effects are presented in the experiment section.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: In the experiment section of the main texts, we discussed several possible countermeasures against our poisoning attack.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We properly credit the origins of open-source datasets, T2I models, and code of RLHF algorithms in the experiment section and appendix.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLM is not used as a necessary core component in this paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.