
TRAP: Hijacking VLA CoT-Reasoning via Adversarial Patches

Zhengxian Huang¹ Wenjun Zhu¹ Haoxuan Qiu² Xiaoyu Ji¹✉ Wenyuan Xu¹

Abstract

By integrating Chain-of-Thought (CoT) reasoning, Vision-Language-Action (VLA) models have demonstrated strong capabilities in robotic manipulation, particularly by improving generalization and interpretability. However, the security of CoT-based reasoning mechanisms remains largely unexplored. In this paper, we show that CoT reasoning introduces a novel attack vector for targeted behavior hijacking—for example, causing a robot to mistakenly deliver a knife to a person instead of an apple—without modifying the user’s instruction. We first provide empirical evidence that CoT strongly governs action generation, even when it is semantically misaligned with the input instructions. Building on this observation, we propose **TRAP**, the first targeted behavior-hijacking adversarial attack against CoT-reasoning VLA models. By targeting the reasoning-to-action pathway, **TRAP** uses an adversarial patch (*e.g.*, a tablecloth placed on the table) to steer intermediate CoT reasoning and downstream actions toward adversary-defined behaviors. Extensive evaluations on three representative reasoning VLAs, spanning distinct CoT reasoning mechanisms, demonstrate the effectiveness of **TRAP**. Notably, we implemented the patch by printing it on paper in a real-world setting. Our findings highlight the urgent need to secure CoT reasoning in VLA systems. The project page is available at [TRAP-website](#).

1. Introduction

Vision–Language–Action (VLA) models (Kim et al., 2024) have reshaped robot learning by enabling open-world manipulation through end-to-end training on large-scale robot data. Recently, a growing body of work has explored rea-

soning VLAs (Zawalski et al., 2024; Black et al., 2025), which incorporate Chain-of-Thought (CoT) reasoning to improve generalization to new scenes and tasks. Beyond performance, explicit reasoning capabilities are also argued to enhance interpretability and safety (Zawalski et al., 2024) by involving intermediate decision-making steps.

Although CoT-enabled reasoning VLA models have attracted growing attention, their security remains largely unexplored. Existing studies primarily examine adversarial attacks on vanilla VLA systems by perturbing perception (Lu et al., 2025b; Yan et al., 2025) or disrupting action generation (Wang et al., 2025a; Jones et al., 2025). Importantly, most prior attacks are untargeted, leading to generic performance degradation. Different from existing attack vectors, CoT often makes the model’s task intent explicit, effectively leaking high-level goals and intermediate plans. In manipulation tasks, CoT can encode intermediate intent as target bounding boxes, motion sequences, or trajectories. This motivates a central question of this work: Can an attacker hijack a VLA model’s behavior by manipulating its CoT reasoning?

In this paper, we show that CoT reasoning can indeed be used for targeted behavior hijacking—for example, causing a robot to mistakenly deliver a knife to a person instead of an apple—without modifying the user’s instruction, as illustrated in the threat scenario of Fig. 2. To investigate how CoT reasoning influences the action generation of VLA models, we conducted a preliminary experiment by actively intervening in instruction–CoT pairs. Our empirical results reveal that CoT strongly governs action generation, even when it is semantically misaligned with the input instructions. Motivated by this insight, we propose **TRAP** (CoT-Reasoning Adversarial Patch), the first targeted behavior-hijacking adversarial attack for reasoning VLAs. **TRAP** uses an adversarial patch (*e.g.*, a tablecloth placed on the table) to corrupt intermediate CoT reasoning and hijack the VLA’s output. Unlike attacks based solely on action losses, **TRAP** induces specific and adversary-defined behaviors by optimizing the CoT adversarial loss. We validate **TRAP** on three VLA architectures, each adopting a distinct CoT reasoning paradigm and collectively covering both integrated and hierarchical CoT mechanisms, demonstrating the broad applicability of our attack. Notably, we demonstrate real-world attacks with visually plausible printed patches,

¹Zhejiang University, Hangzhou, China ²Harbin Institute of Technology, Harbin, China. Correspondence to: Xiaoyu Ji <xji@zju.edu.cn>.

Proceedings of the 43rd International Conference on Machine Learning, Seoul, South Korea. PMLR 306, 2026. Copyright 2026 by the author(s).

enabled by stealthiness-aware optimization, homography transformation, and color calibration. Our findings highlight the urgent need to secure CoT reasoning in VLA systems.

Our main contributions are summarized as follows:

- **New attack vector discovery.** We identify a novel attack vector in CoT reasoning VLAs, where the introduction of CoT mechanisms expands the attack surface and makes it possible to manipulate the behavior of VLAs via adversarial patches.
- **Targeted behavior-hijacking attack against reasoning VLAs.** We propose TRAP, a targeted adversarial attack framework that optimizes a physical patch to hijack intermediate CoT reasoning and steer VLA actions toward attacker-specified behaviors without modifying the user’s instruction. To the best of our knowledge, this is the first adversarial attack that exploits CoT reasoning to induce targeted behavior hijacking in VLAs.
- **Extensive evaluation across reasoning VLA paradigms.** We evaluate TRAP on three representative reasoning VLAs, spanning diverse CoT reasoning mechanisms and model architectures. Results from simulation and real-world physical deployments demonstrate the effectiveness and broad applicability of our attack.

2. Background and Related Work

2.1. Reasoning Vision-Language-Action Models

Vision-Language-Action Models (VLAs) (Zitkovich et al., 2023; Kim et al., 2024) are multi-modal models that can process visual and instruction inputs and generate robot actions. Trained on a large-scale robotic dataset (O’Neill et al., 2024), VLAs show promising performance on general daily-life tasks (Black et al., 2024). Despite the significant progress in VLAs, they still struggle with limited generalization as well as limited performance on long-horizon tasks.

CoT-reasoning VLAs. Inspired by the success of Chain-of-Thought (CoT) reasoning in LLMs, recent VLA models have incorporated embodied CoT mechanisms, a paradigm often referred to as “reasoning VLAs”. Specifically, after perceiving visual and instruction inputs, the reasoning VLA first generates a CoT of diverse forms, such as (1) subtask decomposition (Black et al., 2024; Chen et al., 2025b; Yang et al., 2025), (2) object bounding boxes (Deng et al., 2025), or (3) predicted traces (Lee et al., 2025). Subsequently, the final action is generated by conditioning both the CoT and the inputs. Thus, the generation of CoT-reasoning VLAs can be formulated as:

$$r_t \sim P_\theta(r_t | r_{<t}, O, I) \tag{1}$$

$$a \sim P_\theta(a | R, O, I) \tag{2}$$

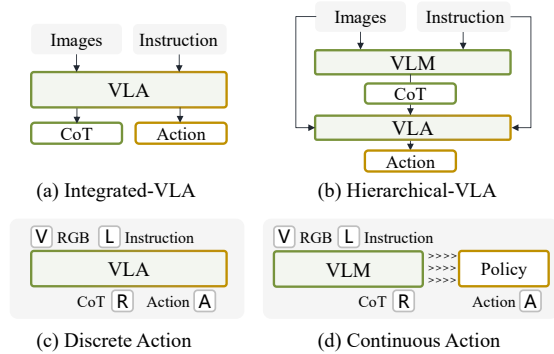


Figure 1. Illustration of existing representative paradigms of reasoning VLAs.

where P_θ denotes the reasoning VLA parameterized by θ . O represents the environmental observations (including visual inputs), and I is the user instruction. Furthermore, r_t denotes the t -th token of the generated CoT R , while a corresponds to the predicted robot action.

CoT-reasoning VLAs can be broadly categorized into Integrated-VLAs and Hierarchical-VLAs (Gao et al., 2026), as illustrated in Figure 1. The Integrated-VLA is a unified model that jointly performs reasoning and action generation. It can be further divided into discrete-token methods (Zawalski et al., 2024; Sun et al., 2025), such as MolmoAct (Lee et al., 2025), which quantize actions into discrete bins and model them as tokens via autoregressive generation, and continuous-regression methods (Black et al., 2025; Yin et al., 2025), such as GraspVLA (Deng et al., 2025), which use diffusion models, flow matching, or MLP heads to directly regress continuous actions. In contrast, the Hierarchical-VLA (Chen et al., 2025b; Shi et al., 2025), such as InstructVLA (Yang et al., 2025), typically adopts a dual-system design: a VLM serves as a high-level planner that generates CoT, which then guides a low-level policy (e.g., a vanilla VLA) to execute precise actions.

CoT mechanisms can improve the performance of VLAs, but they also introduce a new attack vector. We present the first evidence that an attacker can hijack the CoT reasoning process to manipulate a VLA’s downstream actions and induce targeted behaviors. We further show that this vulnerability is general, persisting across diverse CoT paradigms and VLA architectures, including both discrete and continuous action frameworks.

2.2. Adversarial Attacks against VLAs

Adversarial attacks against VLAs have attracted significant attention, particularly because compromised VLAs may cause catastrophic physical consequences in real-world environments. Prior work has primarily focused on **untargeted attacks**, which introduce adversarial perturbations into visual or language inputs to induce task failures. Ex-

isting untargeted attacks include (1) perception-based approaches (Lu et al., 2025b; Yan et al., 2025; Xu et al., 2025a), which undermine scene understanding and cross-modal alignment, and (2) action-based approaches (Wang et al., 2025a; Jones et al., 2025; Wang et al., 2025b), which directly disrupt action generation. For example, RoboticAttack (Wang et al., 2025a) optimizes adversarial patches using action-guided loss objectives, whereas UPA-RFAS (Lu et al., 2025b) attacks visual-encoder representations via a unified feature-space objective coupled with attention steering.

However, targeted behavior-hijacking adversarial attacks against VLAs remain largely underexplored, despite posing a substantially greater threat by inducing targeted malicious behaviors. Moreover, existing studies have primarily focused on vanilla VLA architectures, leaving the vulnerabilities of emerging CoT-reasoning VLAs largely overlooked. In this paper, we bridge this gap by analyzing the security of CoT-reasoning VLAs and designing targeted attacks that explicitly manipulate VLAs’ behavior.

3. Threat Model

Attack Goal. The attacker’s primary objective is to induce a VLA to execute a *targeted* malicious behavior (e.g., handing a knife to the user), as illustrated in Figure 2. Concretely, the attacker can introduce a printed adversarial patch into the scene (e.g., a tablecloth on the table), which biases the model’s CoT reasoning and consequently causes malicious or unintended behavior. The attack is further designed to persist across multiple VLA inference steps, reflecting the short-horizon nature of current VLAs, which typically predict only one or a few actions per inference.

Adversary Capability. Following prior work (Wang et al., 2025a), we consider a white-box adversary who has access to the victim VLA’s architecture, parameters, and gradients. In the black-box setting, the attacker can transfer an adversarial patch optimized on a surrogate VLA. We further assume the adversary can physically deploy the patch in the environment (e.g., by sticking the patch on the table). For reliable physical realization, the adversary may obtain deployment-specific calibration priors, such as homography estimation and a color calibration model, but these priors do not grant any control beyond physical patch placement. Importantly, the adversary cannot manipulate the instruction input, i.e., all instructions remain benign user-provided commands.

4. Preliminary Analysis

Despite the widespread adoption of CoT in recent VLA models (Deng et al., 2025; Chen et al., 2025b; Black et al., 2025), its underlying mechanisms, particularly its security implica-

tions, remain underexplored. ECoT-lite (Chen et al., 2025a) suggests that ECoT’s gains primarily stem from improved representation learning induced by reasoning supervision, and VLA-OS (Gao et al., 2026) provides a comprehensive empirical study of different CoT paradigms and their representational modalities. However, existing work has not disentangled the causal roles of instructions and CoT, either individually or jointly, in driving action generation. More critically, when instructions and CoT conflict, the internal arbitration dynamics remain unclear, leaving open which signal dominates the final action. These gaps motivate two key questions for enabling targeted behavior-hijacking attacks:

- **Role of CoT in action generation:** How much causal influence does CoT exert on downstream action generation, and can manipulating CoT reliably steer the model toward targeted actions?
- **Conflict arbitration:** When input instructions and CoT are misaligned, does the model adhere more faithfully to the original instruction or to the intermediate reasoning?

To answer these two questions, we conduct a preliminary analysis of the causal mechanisms in reasoning VLAs. Specifically, we perform active interventions on the inputs of the VLAs and observe the subsequent action generation. We first collect original tuples (I, R) and then consider two experimental settings:

(1) Instruction Masking. We assess the impact of CoT R on action generation by masking instruction tokens, $P_{\theta}(a|R, O, \cdot)$. This setting enables us to examine whether VLAs can complete tasks when provided with only CoT.

(2) Cross-Sample Shuffling. To create semantic conflict, we disrupt the correspondence between the instruction and the CoT. For the i -th tuple, we replace its instruction $I^{(i)}$ with a distractor $I^{(j)}$ from a different task instance $j \neq i$ (e.g., $I^{(i)}$ is “pick apple”, and its corresponding CoT R^i is the intermediate reasoning information for reaching the apple, while $I^{(j)}$ is “pick knife”). The action is generated via $P_{\theta}(a|R^{(i)}, O, I^{(j)})$. This setting creates conflict where $R^{(i)}$ and $I^{(j)}$ suggest different follow-up action patterns, allowing us to investigate which one dominates, and it is closer to the settings of our subsequent attack.

Results and Analysis. We examine three mainstream reasoning VLAs across various tasks to evaluate how their behaviors change under the two intervention settings relative to the original configuration (See Section 6.1 for full details on the experimental environment and metrics.). As shown in Table 1, we observe the task performance degradation under the two aforementioned settings compared to TSR_{ori} , suggesting that instructions and CoT jointly con-

Table 1. Instruction masking and CoT shuffling analysis on reasoning VLAs. TSR denotes Task Success Rate. Subscripts m, s^i, s^j denote mask and shuffling settings respectively.

Model	TSR _{ori}	TSR _m	Sim _m	TSR _{sⁱ}	TSR _{s^j}	Score
InstructVLA	45.23%	39.95%	0.3169	16.42%	13.43%	0.0508 ± 0.2392
MolmoAct	50.40%	17.35%	0.1805	12.94%	16.26%	0.0520 ± 0.1964
GraspVLA	94.20%	92.60%	0.3990	94.20%	0.00%	0.5278 ± 0.1696

tribute positively to action generation. Under the instruction masking setting, the VLAs retain partial capability when driven solely by CoT, highlighting CoT’s positive influence on action generation. Under the cross-sample shuffling setting, TSR_{sⁱ} shows that CoT plays a more dominant role than the instruction, whereas TSR_{s^j} suggests the opposite. Notably, the near-zero scores for both InstructVLA and MolmoAct imply that, across tasks, instruction and CoT act as competing signals with roughly comparable influence. In contrast, the high score achieved by GraspVLA suggests that, for this model, CoT plays a dominant role in governing action generation. Additional experimental details and analysis are provided in Section B.2.

Key Observation: *CoT-reasoning VLAs exhibit a “competition mechanism” in which input instructions and CoT competitively influence action generation, with CoT acting as a vital driver of the final action.*

5. Methodology

Motivated by the pivotal role of CoT in reasoning VLAs in Section 4, we propose an adversarial patch attack designed to manipulate VLAs’ behavior via CoT hijacking. In the following sections, we formalize the problem statement and detail our proposed attack.

5.1. Problem Formulation

Adversarial Patch Injection. The adversary aims to manipulate the VLAs to execute a target behavior by applying an adversarial patch to the environment. Let δ denote the adversarial patch and M be a binary mask indicating the patch’s position. For every time step t , the adversarial observation \tilde{O} is formulated as:

$$\tilde{O} = (1 - M) \odot O + M \odot \delta \quad (3)$$

where \odot denotes element-wise multiplication. The adversarial patch should remain effective throughout the entire rollout. **Offline Rollout.** Unlike previous patch attack work on a single static frame, our goal is to hijack the VLAs to execute a coherent, sequential target behavior. Thus, we curate an offline clean dataset $\mathcal{D} = \{\tau_n\}_{n=1}^N$ by rolling out the VLAs on different tasks. The trajectory τ is represented as a tuple (O, R, a) .

Optimization Objective. The adversarial objective is to

optimize the patch δ such that when it is applied to the clean observations in \mathcal{D} , the VLA P_θ is induced to generate the target reasoning CoT R^* and action a^* . Additionally, to improve the stealthiness of the patch, we further introduce a content loss $\mathcal{L}_{\text{content}}$ and Total Variation (TV) loss \mathcal{L}_{TV} . The coefficients λ_1, λ_2 , and λ_3 balance the contributions of the action loss, content loss, and TV loss, respectively. Formally, we minimize:

$$\min_{\delta \in \Delta} \mathbb{E}_{\tau \sim \mathcal{D}} [\mathcal{L}_{\text{cot}} + \lambda_1 \mathcal{L}_{\text{action}} + \lambda_2 \mathcal{L}_{\text{content}} + \lambda_3 \mathcal{L}_{\text{TV}}] \quad (4)$$

5.2. Attack Effectiveness

CoT Hijacking Loss. Prevailing reasoning VLAs predominantly employ VLMs to synthesize CoT reasoning. Despite variations in the specific form of CoT, its generation is commonly formulated as a standard next-token prediction task. Thus, we employ a Cross-Entropy (CE) loss to align the generated CoT tokens with the target sequence R^* .

$$\mathcal{L}_{\text{cot}} = - \sum_{t=1}^T \log P_\theta(r_t^* | r_{<t}^*, \tilde{O}, I) \quad (5)$$

Action Loss. As discussed in Section 2.1, existing VLA architectures utilize a variety of action-decoding heads—typically discrete tokenization or continuous regression. Therefore, we apply different action loss strategies. For VLAs that treat actions as discrete tokens, we utilize the CE loss:

$$\mathcal{L}_{\text{action}}^{\text{disc}} = - \log P_\theta(a^* | R^*, \tilde{O}, I) \quad (6)$$

For VLAs that regress to continuous action (e.g., diffusion-based), especially those employing action chunking (Zhao et al., 2023), we transform the action sequence into waypoints. We employ the Mean Squared Error (MSE):

$$\mathcal{L}_{\text{action}}^{\text{cont}} = \|f_{\text{traj}}(a) - f_{\text{traj}}(a^*)\|_2^2 \quad (7)$$

where $f_{\text{traj}}(\cdot)$ maps the raw model output to the trajectory.

To optimize the adversarial patch δ , we adopt Projected Gradient Descent (PGD) (Madry et al., 2017) and the optimization process can be formulated as:

$$\delta_{t+1} = \text{Proj}_{\|\cdot\|_\infty \leq \epsilon} (\delta_t + \eta \cdot \nabla_\delta L(\delta_t)) \quad (8)$$

where $\text{Proj}(\cdot)$ projects δ into perturbation budget ϵ , η is the attack step size, and the $\nabla L(\cdot)$ is the gradient of total loss.

5.3. Stealthiness Enhancement

Content Loss. Without camouflage constraints, adversarial patch optimization tends to produce semantically meaningless patterns that are visually conspicuous to humans. To improve the stealthiness of the attack, we follow (Zhu et al., 2023) and introduce a content loss. Specifically, given a

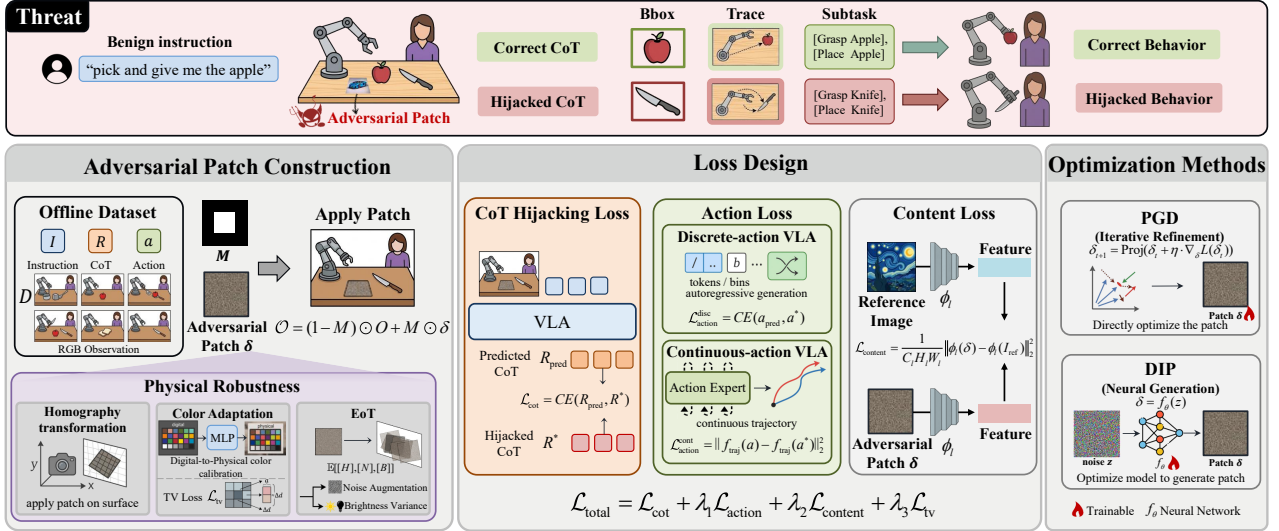


Figure 2. Overview of the TRAP attack framework. The adversary places an adversarial patch in the scene to corrupt the VLA’s intermediate CoT, causing the model to execute an attacker-specified behavior while the user instruction remains benign. The patch is optimized with two groups of objectives: attack-effectiveness losses, including CoT hijacking and action losses, and stealthiness losses, including content and TV losses.

reference image I_{ref} (e.g., an image of a sports car) and a patch δ of the same size $C \times H \times W$, we use the l -th layer of a pretrained CNN, denoted by ϕ_l , to extract their feature representations. We then apply the MSE loss to encourage the patch to capture the content and spatial structure of the reference image.

$$\mathcal{L}_{\text{content}} = \frac{1}{C_l H_l W_l} \|\phi_l(\delta) - \phi_l(I_{\text{ref}})\|_2^2 \quad (9)$$

TV Loss. We employ a Total Variation (TV) loss to penalize high-frequency artifacts and encourage color consistency, thereby ensuring the physical realizability of the adversarial patch as well as stealthiness.

$$\mathcal{L}_{\text{tv}} = \sum_{i,j} \sqrt{(x_{i,j} - x_{i+1,j})^2 + (x_{i,j} - x_{i,j+1})^2} \quad (10)$$

where $x_{i,j}$ is the pixel value of the adversarial patch δ at the coordinate (i, j) .

DIP Optimization. Inspired by Deep Image Prior (DIP) (Ulyanov et al., 2018), we leverage the implicit regularization of a CNN to synthesize adversarial patches that exhibit spatially coherent and visually less high-frequency noise. Specifically, instead of directly optimizing the patch δ in pixel space, we optimize the parameters θ of a CNN f_θ that maps a fixed noise input z to the adversarial patch, i.e., $\delta = f_\theta(z)$.

5.4. Physical Robustness

Homography Transformation. To guarantee that the adversarial patch is compatible with the physical world rather than directly applied on the image plane, we model the ge-

ometric transformation of the adversarial patch from the table plane to the image plane as a planar homography. Let $\mathbf{p} = [u, v, 1]^T$ denote the homogeneous coordinates of a pixel in the canonical patch space, and $\mathbf{p}' = [x, y, 1]^T$ be the corresponding coordinates in the target image space. The mapping is defined up to a scale factor by a 3×3 matrix \mathbf{H} :

$$\mathbf{p}' \sim \mathbf{H}\mathbf{p} \iff \lambda \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \quad (11)$$

Color Calibration. Color distortion inherent in the print-and-capture process often degrades the efficacy of physical adversarial patches. To mitigate this, we employ a Multi-Layer Perceptron (MLP) to model the transformation from digital simulation to the physical domain. During the optimization phase, this model serves as a calibration function, aligning the digital patch’s color distribution with physical reality.

Expectation over Transformation. To improve the robustness of the patch in the physical world, we adopt the Expectation over Transformation (EoT) method (Athalye et al., 2018), which optimizes the patch over a distribution of transformations.

6. Experiment

6.1. Experimental Setup

VLA Models. We evaluate three representative reasoning VLAs: MolmoAct (Lee et al., 2025), GraspVLA (Deng et al., 2025), and InstructVLA (Yang et al., 2025). These

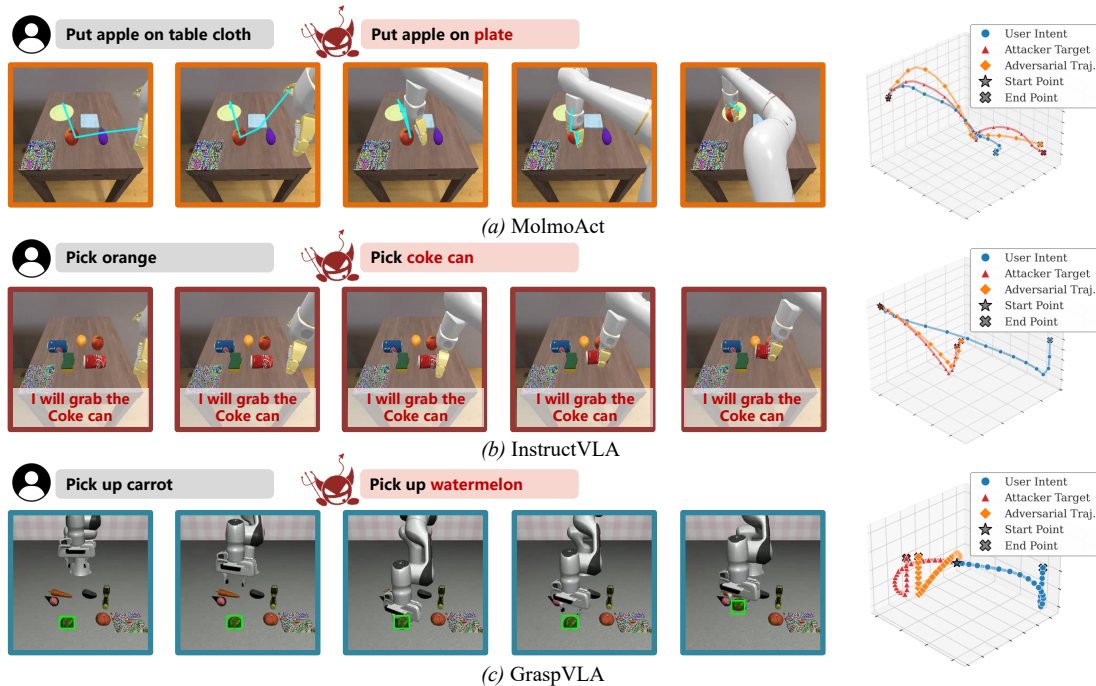


Figure 3. **Qualitative results of TRAP attack.** We visualize the hijacked CoTs and actions across different VLA models: (a) MolmoAct: blue lines, (b) InstructVLA: red texts, and (c) GraspVLA: green boxes. The visualized 3D trajectories demonstrate that **TRAP** effectively hijacks various VLAs to execute the attacker’s target behaviors.

Table 2. **Description of victim VLAs.** “Disc.” denotes Discrete, “Cont.” denotes Continuous.

Model	Arch.	CoT Type
MolmoAct	Integrated (Disc.)	Trace, Depth
GraspVLA	Integrated (Cont.)	Bbox, Grasp pose
InstructVLA	Hierarchical (Cont.)	Subtasks

models cover diverse architectural designs and CoT reasoning paradigms, as summarized in Table 2. Detailed descriptions are provided in Appendix C.1.

Tasks. We construct a suite of manipulation-task pairs for evaluation. For each task, we generate two instructions: one representing the user’s intended task and the other serving as the attack target. The former is consistently provided to the victim VLAs and remains unchanged; meanwhile, the adversary seeks to place an adversarial patch to hijack the VLAs to perform the behavior specified by the latter. Details are provided in Appendix C.2.

Evaluation Metrics. We define the Attack Success Rate (ASR) as the completion rate of target tasks. Task completion is checked via physical contact by the simulators (Li et al., 2024; Liu et al., 2023). A trial is counted as successful only when the robot completes the attacker-specified target behavior. To complement the binary task-completion signal, we propose a fine-grained score to quantify the extent of hijacking for a trajectory:

$$\text{Score}(T) = \frac{\text{sim}(T, T_A) - \text{sim}(T, T_B)}{\text{sim}(T, T_A) + \text{sim}(T, T_B)} \quad (12)$$

where T is the trajectory to evaluate, T_A is the target trajectory, T_B is the original trajectory of the user instruction, $\text{sim} \in [0, 1]$ is the similarity between two trajectories, which is implemented with Dynamic Time Warping (DTW).

Compared Methods. We compare TRAP with two representative baselines: (a) Random Noise, serving as a lower bound, and (b) Action Attack, which represents a standard end-to-end adversarial attack against the VLA, e.g., TMA (Wang et al., 2025a). In addition, we report TRAP_{CoT-only}, a reasoning-only variant of our method that optimizes only the proposed CoT hijacking loss, i.e., without the explicit action loss.

Evaluation Scope. Our simulation study focuses on attack effectiveness rather than visual stealthiness. This design allows us to isolate the core question of whether reasoning manipulation can induce targeted VLA behaviors. Accordingly, we report ASR and trajectory hijack score as the primary metrics, while leaving stealthiness-oriented objectives and visual camouflage evaluation to the physical-world setting.

Implementation Details. For the dataset and simulation evaluation, following SimplerEnv (Li et al., 2024), we sample 25 different object layouts per task instance for training, and 10 more layouts for testing. Each task is executed for 5 trials, resulting in a total of 175 rollouts per task. In addition, the maximum number of epochs is set to 80 in SimplerEnv, whereas it is 100 in LIBERO. For the adversarial patch opti-

mization, the pixel update step is set to 8/255, with a batch size of 4. For simulation evaluation, we adopt different loss weights λ_1 for VLAs, and λ_1 is set to 1 for MolmoAct and GraspVLA and to 2 for InstructVLA. For real-world evaluation, we set λ_2 to 1.5 and λ_3 to 100 for stealthiness. For DIP optimization, we adopt an annealed regularization strategy that applies stronger visual regularization in the early stage for stealthiness and gradually decays these weights in later iterations to emphasize the VLA attack objective. All the adversarial patch generation experiments are conducted on one NVIDIA H800 GPU, while the simulation evaluation experiments run with one RTX 4090 GPU.

6.2. Overall Performance

Main Results. We optimize one adversarial patch per task, leveraging the data from 25 different object layout instances and further evaluate it on 10 unseen layout instances. As shown in Table 3, **TRAP** consistently outperforms all baselines on InstructVLA and GraspVLA by a large margin. On MolmoAct, **TRAP** achieves performance comparable to that of TRAP_{CoT-only} (48.06% vs. 49.52%). However, we observe that **TRAP** demonstrates higher hijack scores (0.3904 vs. 0.3422), suggesting that its generated trajectories better conform to the VLA’s nominal behavior patterns.

Notably, on the MolmoAct and GraspVLA, both **TRAP** and the TRAP_{CoT-only} significantly surpass the Action Attack method. This underscores the critical role of CoT in reasoning VLAs: once the CoT is successfully hijacked, the subsequent action predictions become highly susceptible to manipulation. However, the efficacy of the TRAP_{CoT-only} approach diminishes on InstructVLA. A closer inspection shows that, although InstructVLA can generate the target CoT, *i.e.*, the desired subtask predictions, its action outputs exhibit severe mode collapse and often degenerate into repetitive behaviors. In contrast, **TRAP** mitigates this issue by incorporating an explicit action loss, which constrains the model toward a plausible behavioral manifold. These results suggest that TRAP_{CoT-only} is effective when corrupted CoT can be reliably propagated to the action decoder, while **TRAP** is necessary when the CoT-to-action coupling is weak or unstable, where action-level supervision helps translate hijacked reasoning into executable target behaviors.

Crucially, **TRAP** exhibits remarkable generalization capabilities. When evaluated on 10 unseen object layouts, **TRAP** maintains near-identical performance to the training instances (*e.g.*, 51.60% vs. 52.54% average ASR). This minimal performance decay demonstrates that **TRAP** learns layout-invariant adversarial features rather than over-fitting to specific spatial configurations.

Interpretability Analysis. As illustrated in Fig. 4, we visualize the attentional allocation of the VLA over image tokens during the generation of CoT. By comparing Fig. 4b and

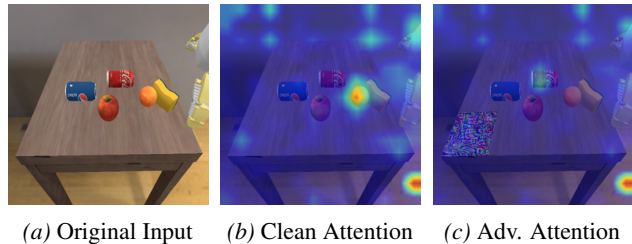


Figure 4. **Illustration of the attention shift.** The 19th layer attention of MolmoAct is visualized. The patch is optimized to redirect the model’s execution from the benign instruction “pick orange” to the adversarial target “pick coke can”.

Fig. 4c, we observe a conspicuous transition in the attention distribution, where the model’s focus is redirected from the benign user intent (“pick orange”) toward the adversarial target (“pick coke”). This attention shift remains consistent across varying spatial configurations of the objects and VLAs, indicating that the attack successfully establishes a spurious mapping between semantic concepts and visual features.

6.3. Further Analysis

Impact of Instruction Variants. In a realistic scenario, the user may use different instructions to denote the same intent. Therefore, we generate instruction variants in two ways: (i) **Paraphrasing**, where we perturb the syntactic structure of the instruction (*e.g.*, “grasp the orange” instead of “pick orange”); and (ii) **Extra Context**, where the core instruction is augmented with irrelevant or environmental descriptions (*e.g.*, “go ahead and pick orange” instead of “pick orange”). The details of the instruction variant are provided in Table 9. Due to the limited variance in GraspVLA’s instruction set (predominantly restricted to “pick up” objects), we omit it from this experiment and prioritize the evaluation of InstructVLA and MolmoAct. Meanwhile, we select the task with the highest attack success rate across the evaluated models to serve as the primary task in this experiment.

As reported in Table 4, **TRAP** remains effective across instruction variants. Notably, it achieves ASR of 70.6% and 60.0% on MolmoAct, while consistently maintaining a score exceeding 0.3. These findings suggest that optimization builds a context-dependent association between the adversarial patch and target-object names, rather than a fixed instruction-template trigger. In matching scenes, these object-level cues can still bias the model toward the adversarial target, increasing practical risk under natural instructions. Compared with MolmoAct, InstructVLA exhibits weaker yet still effective transferability. We hypothesize that InstructVLA is more sensitive to instruction variants because its CoT relies on text subtask decomposition, whereas MolmoAct is more robust to linguistic variation due to its reliance on trajectory-related tokens for internal reasoning.

Table 3. **Main results:** Overall performance of TRAP attacks on reasoning VLAs across all tasks. TRAP_{CoT-only} denotes our reasoning-only variant that optimizes only the CoT hijacking loss, while TRAP jointly optimizes CoT and action objectives.

Method	MolmoAct		InstructVLA		GraspVLA		Average	
	ASR(%)	Score	ASR(%)	Score	ASR(%)	Score	ASR(%)	Score
Random Noise	0.97	-0.3769	3.39	-0.3280	0.32	-0.3064	1.56	-0.3371
Action Attack	9.68	0.1282	6.77	-0.2736	0.00	-0.2948	5.48	-0.1467
TRAP _{CoT-only}	49.52	0.3422	4.03	-0.0332	69.04	0.3898	40.86	0.2329
TRAP	48.06	0.3904	33.71	0.1724	75.84	0.4253	52.54	0.3294
↳Unseen layouts	48.00	0.1828	31.60	0.1307	75.20	0.4024	51.60	0.2386

Table 4. **Instruction analysis:** Performance under Paraphrasing and Extra-Context instruction variations.

	MolmoAct		InstructVLA	
	ASR(%)	Score	ASR(%)	Score
Original	72.0	0.34	67.4	0.56
Paraphrasing	70.6	0.32	25.1	0.09
Extra-Context	60.0	0.35	44.8	0.32

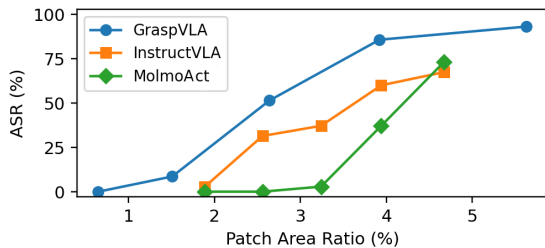


Figure 5. Impact of patch size on attack effectiveness.

Impact of Patch Size. We evaluate TRAP’s performance with different patch sizes in Fig. 5. Empirical results demonstrate that performance metrics scale positively with increasing patch size. It is reasonable that the larger patch size provides a larger optimization space, aligning with prior work. Notably, we observe a behavioral shift in MolmoAct as the patch size increases. At around 1.9%, MolmoAct largely follows the user intent, whereas at 2.6% it reaches an apparent inflection point where the robot remains nearly static. Beyond this point, the model’s behavior shifts more decisively toward the attacker’s objective.

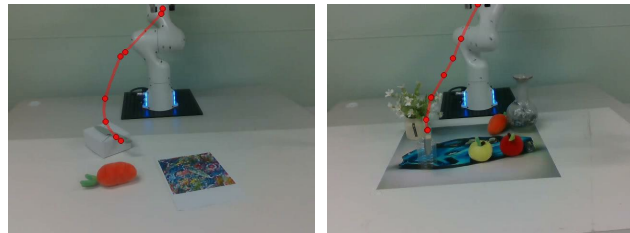
6.4. Transferability Study

We further investigate the transferability of our proposed attack across different fine-tuning and input-state configurations within the VLA model families. For MolmoAct, we evaluate cross-checkpoint transfer by optimizing the attack on the RT-1-fine-tuned checkpoint and deploying the same attack against the pretrained, non-RT-1-fine-tuned checkpoint. For InstructVLA, we study transfer across robot-state configurations by optimizing the attack on the variant without robot-state input and evaluating it on the corresponding state-conditioned variant.

As shown in Table 5, TRAP exhibits a certain degree of transferability. This transferability is particularly concern-

Table 5. **Transferability results.** Performance comparison between source models (optimized) and target variants (transfer).

Model Pair	Variant	ASR(%) ↑	Score ↑
MolmoAct	Optimized	48.05	0.329
	Transfer	18.39	0.170
InstructVLA	Optimized	33.10	0.160
	Transfer	9.20	-0.170



(a) Occlusion-free deployment (b) Object-occluded deployment

Figure 6. **Physical deployment settings for real-world evaluation.** We evaluate (a) an occlusion-free tabletop patch and (b) an object-occluded tablecloth patch. The benign instruction is “pick up carrot”, while the target is “pick up knife”; red waypoints denote the hijacked trajectory. Videos of the real-world experiments are available on our [project homepage](#).

ing under the prevailing deployment paradigm. Since pre-training VLA models is computationally expensive, practitioners commonly adapt publicly released pretrained checkpoints to private domains through downstream fine-tuning. Therefore, an attacker can conduct an adversarial attack on a publicly available surrogate checkpoint that shares the same architecture or initialization, and then transfer it to a privately fine-tuned model in a black-box setting.

6.5. Real-world Evaluation

For real-world evaluation, we conduct the real-world evaluation on GraspVLA, following its official physical setup. We consider two physical experimental settings, as illustrated in Fig. 6: (1) an occlusion-free deployment and (2) an object-occluded deployment. We designed a “hazardous redirection” scenario where an adversary attempts to hijack the user’s original intent (“pick up carrot”) toward a malicious target (“pick up knife”).

Occlusion-free Deployment. In this setting, the optimized adversarial patch is placed flat on the tabletop, where it

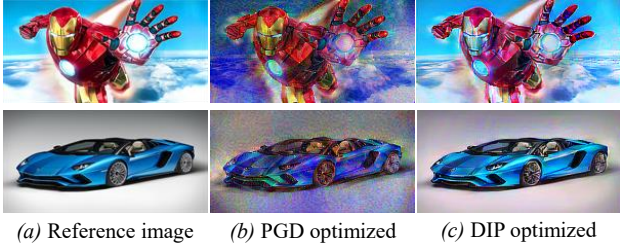


Figure 7. Impact of optimization methods on real-world adversarial patches.

remains unoccluded by surrounding objects throughout execution. In an evaluation spanning 15 independent trials, the generated physical patch demonstrated significant capability in hijacking the VLA model’s execution. The attack successfully subverted the model’s internal reasoning for at least one reasoning step in **13/15 (86.7%)** of cases. In **5/15 (33.3%)** of trials, the adversary maintained control over the entire physical manipulation process, successfully completing the malicious objective.

Object-occluded Deployment. In this setting, we consider a more realistic deployment, where the adversarial patch is used as a tablecloth or placemat, and objects are placed on top of it, naturally inducing object clutter and partial occlusion. To enhance stealthiness, the patch is rendered as a semantically meaningful pattern using the methods described in Section 5.3.

We compare adversarial patches optimized with PGD and DIP in terms of both visual stealthiness and real-world attack effectiveness. As shown in Fig. 7, DIP produces smoother and more spatially coherent adversarial patterns than direct pixel-space optimization with PGD, better preserving the semantic appearance of the reference image due to the implicit regularization of CNNs. We further evaluate both patches under the real-world object-occluded deployment over 50 trials. The PGD-optimized patch achieves an ASR of **19/50 (38.0%)**, while the DIP-optimized patch achieves an ASR of **17/50 (34.0%)**. This suggests that DIP substantially improves visual stealthiness while largely preserving attack effectiveness. We also verify the robustness of the optimized patches under diverse object-layout variations, where the attacks remain effective when objects are added, removed, or repositioned on the patch. Additional visualizations are provided in Appendices D.4 and D.5.

7. Discussion

Countermeasure. TRAP exploits the CoT reasoning process to facilitate adversarial attacks, thereby inducing a structural inconsistency between the user’s original instruction and the CoT. Motivated by this observation, we develop an independent and lightweight detector to verify whether the task intention reflected in the CoT remains aligned with the

Table 6. **Lightweight defense performance.** Parentheses indicate changes relative to the GPT-5-based method. BCR: Benign Consistent Rate; ARR: Attack Reject Rate. Parentheses indicate relative changes compared with GPT-5.

CoT	BCR	ARR	Time
Bbox	87.00% (+12.00%)	92.50% (-1.30%)	201.67 ms
Trace	77.01% (-5.79%)	87.42% (+8.02%)	114.35 ms
Subtask	86.07% (-9.43%)	91.94% (-8.06%)	2.62 ms

initial user instruction. Since different reasoning VLAs represent CoT in different formats, we design format-specific lightweight consistency checks. Specifically, (1) For bbox-based CoT (e.g., GraspVLA), we use an open-vocabulary detector to verify the consistency between the detected target-object bounding box and the bounding box predicted by the VLA; (2) for trace-based CoT (e.g., MolmoAct), we check whether the predicted trace is consistent with the detected bounding box of the target object; (3) For text-based CoT (e.g., InstructVLA), we employ a lightweight text encoder to measure the semantic similarity between the original instruction and the generated reasoning subtask. We further compare our lightweight defense with a strong VLM-based baseline, GPT-5, which offers stronger reasoning capability but incurs substantially higher latency and computational cost.

As shown in Table 6, our proposed lightweight defenses achieve performance close to the GPT-5-based method while incurring only millisecond-level overhead. Compared with the runtime of the reasoning VLAs, the overhead of our lightweight defenses is small; in the most efficient case, it is only 0.20% of the model inference time (for InstructVLA, with nearly 1.3s inference time).

Limitation. Current open-source reasoning VLAs still have limited capabilities in genuinely long-horizon manipulation tasks without task-specific fine-tuning. Accordingly, our evaluation focuses on “pick-and-place” primitives, while our method could potentially extend to long-horizon VLA tasks as reasoning VLAs advance. Additionally, improving the transferability of adversarial patches within the VLA family, as well as across structurally different VLAs, remains an important direction for future work.

8. Conclusion

We introduce TRAP, the first targeted behavior-hijacking adversarial attack tailored for reasoning VLAs. In contrast to existing untargeted perturbations, our method exploits the inherent vulnerabilities within the CoT reasoning mechanisms of VLAs. This enables TRAP to induce precise, adversary-defined behaviors. Comprehensive experiments across diverse VLA architectures demonstrate the efficacy of our approach and highlight a critical security bottleneck in current embodied AI systems.

Acknowledgements

We thank the anonymous reviewers and area chair for their valuable comments and constructive suggestions. We are also grateful to Dr. Richeng Jin for helpful discussions and valuable feedback on this work. This work is supported by the National Natural Science Foundation of China (NSFC) Grant 62222114 and 61925109.

Impact Statement

This work reveals a critical security vulnerability in reasoning Vision-Language-Action models, showing that adversarial patches can hijack intermediate reasoning and induce targeted robotic behaviors without modifying user instructions. While such findings are dual-use, our goal is to raise awareness of potential physical safety risks in embodied AI systems and motivate future research on robust reasoning, consistency checking, and secure deployment. We believe that understanding these vulnerabilities is essential for building safer and more trustworthy embodied AI systems.

References

- Athalye, A., Engstrom, L., Ilyas, A., and Kwok, K. Synthesizing robust adversarial examples. In *International conference on machine learning*, pp. 284–293. PMLR, 2018.
- Black, K., Brown, N., Driess, D., Esmail, A., Equi, M., Finn, C., Fusai, N., Groom, L., Hausman, K., Ichter, B., et al. π_0 : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- Black, K., Brown, N., Darpinian, J., Dhabalia, K., Driess, D., Esmail, A., Equi, M. R., Finn, C., Fusai, N., Galliker, M. Y., et al. $\pi_{0.5}$: a vision-language-action model with open-world generalization. In *9th Annual Conference on Robot Learning*, 2025.
- Chen, W., Belkhale, S., Mirchandani, S., Mees, O., Driess, D., Pertsch, K., and Levine, S. Training strategies for efficient embodied reasoning. *arXiv preprint arXiv:2505.08243*, 2025a.
- Chen, X., Chen, Y., Fu, Y., Gao, N., Jia, J., Jin, W., Li, H., Mu, Y., Pang, J., Qiao, Y., et al. Internvla-m1: A spatially guided vision-language-action framework for generalist robot policy. *arXiv preprint arXiv:2510.13778*, 2025b.
- Deng, S., Yan, M., Wei, S., Ma, H., Yang, Y., Chen, J., Zhang, Z., Yang, T., Zhang, X., Zhang, W., et al. Graspvla: a grasping foundation model pre-trained on billion-scale synthetic action data. *arXiv preprint arXiv:2505.03233*, 2025.
- Gao, C., Liu, Z., Chi, Z., Huang, J., Fei, X., Hou, Y., Zhang, Y., Lin, Y., Fang, Z., and Shao, L. Vla-os: Structuring and dissecting planning representations and paradigms in vision-language-action models. *Advances in Neural Information Processing Systems*, 38:136705–136736, 2026.
- Guo, J., Jiang, W., Lin, Y., Liu, Y., Zhang, R., Lu, G., Chen, A., Han, X., Li, H., and Niyato, D. State backdoor: Towards stealthy real-world poisoning attack on vision-language-action model in state space. *arXiv preprint arXiv:2601.04266*, 2026.
- Jones, E. K., Robey, A., Zou, A., Ravichandran, Z., Pappas, G. J., Hassani, H., Fredrikson, M., and Kolter, J. Z. Adversarial attacks on robotic vision language action models. *arXiv preprint arXiv:2506.03350*, 2025.
- Kim, M. J., Pertsch, K., Karamcheti, S., Xiao, T., Balakrishna, A., Nair, S., Rafailov, R., Foster, E., Lam, G., Sankeeti, P., et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- Kuo, M., Zhang, J., Ding, A., Wang, Q., DiValentin, L., Bao, Y., Wei, W., Li, H., and Chen, Y. H-cot: Hijacking the chain-of-thought safety reasoning mechanism to jailbreak large reasoning models, including openai o1/o3, deepseek-r1, and gemini 2.0 flash thinking. *arXiv preprint arXiv:2502.12893*, 2025.
- Lee, J., Duan, J., Fang, H., Deng, Y., Liu, S., Li, B., Fang, B., Zhang, J., Wang, Y. R., Lee, S., et al. Molmoact: Action reasoning models that can reason in space. *arXiv preprint arXiv:2508.07917*, 2025.
- Li, J., Zhao, Y., Zheng, X., Xu, Z., Li, Y., Ma, X., and Jiang, Y.-G. Attackvla: Benchmarking adversarial and backdoor attacks on vision-language-action models. *arXiv preprint arXiv:2511.12149*, 2025.
- Li, X., Hsu, K., Gu, J., Pertsch, K., Mees, O., Walke, H. R., Fu, C., Lunawat, I., Sieh, I., Kirmani, S., et al. Evaluating real-world robot manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024.
- Liu, B., Zhu, Y., Gao, C., Feng, Y., Liu, Q., Zhu, Y., and Stone, P. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36:44776–44791, 2023.
- Lu, C., Fan, X., Huang, Y., Xu, R., Li, J., and Xu, W. Does chain-of-thought reasoning really reduce harmfulness from jailbreaking? In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 6523–6546, 2025a.
- Lu, H., Yu, Y., Yang, Y., Yi, C., Zhang, Q., Shen, B., Kot, A. C., and Jiang, X. When robots obey the patch: Universal transferable patch attacks on vision-language-action models. *arXiv preprint arXiv:2511.21192*, 2025b.

- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- O’Neill, A., Rehman, A., Maddukuri, A., Gupta, A., Padalkar, A., Lee, A., Pooley, A., Gupta, A., Mandlekar, A., Jain, A., et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6892–6903. IEEE, 2024.
- Shi, L. X., Ichter, B., Equi, M., Ke, L., Pertsch, K., Vuong, Q., Tanner, J., Walling, A., Wang, H., Fusai, N., et al. Hi robot: Open-ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025.
- Sun, Q., Hong, P., Pala, T. D., Toh, V., Tan, U.-X., Ghosal, D., and Poria, S. Emma-x: An embodied multimodal action model with grounded chain of thought and look-ahead spatial reasoning. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14199–14214, 2025.
- Ulyanov, D., Vedaldi, A., and Lempitsky, V. Deep image prior. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9446–9454, 2018.
- Wang, T., Han, C., Liang, J., Yang, W., Liu, D., Zhang, L. X., Wang, Q., Luo, J., and Tang, R. Exploring the adversarial vulnerabilities of vision-language-action models in robotics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6948–6958, 2025a.
- Wang, X., Li, J., Weng, Z., Wang, Y., Gao, Y., Pang, T., Du, C., Teng, Y., Wang, Y., Wu, Z., et al. Freezevla: Action-freezing attacks against vision-language-action models. *arXiv preprint arXiv:2509.19870*, 2025b.
- Xu, H., Koh, Y. S., Huang, S., Zhou, Z., Wang, D., Sakuma, J., and Zhang, J. Model-agnostic adversarial attack and defense for vision-language-action models. *arXiv preprint arXiv:2510.13237*, 2025a.
- Xu, R., Qi, Z., and Xu, W. Preemptive answer “attacks” on chain-of-thought reasoning. In *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 14708–14726, 2024.
- Xu, Z., Zheng, X., Ma, X., and Jiang, Y.-G. Tabvla: Targeted backdoor attacks on vision-language-action models. *arXiv preprint arXiv:2510.10932*, 2025b.
- Yan, Y., Xie, Y., Zhang, Y., Lyu, L., and Jin, Y. When alignment fails: Multimodal adversarial attacks on vision-language-action models. *arXiv preprint arXiv:2511.16203*, 2025.
- Yang, S., Li, H., Chen, Y., Wang, B., Tian, Y., Wang, T., Wang, H., Zhao, F., Liao, Y., and Pang, J. Instructvla: Vision-language-action instruction tuning from understanding to manipulation. *arXiv preprint arXiv:2507.17520*, 2025.
- Yao, Y., Tong, X., Wang, R., Wang, Y., Li, L., Liu, L., Teng, Y., and Wang, Y. A mousetrap: Fooling large reasoning models for jailbreak with chain of iterative chaos. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 7837–7855, 2025.
- Yin, C., Lin, Y., Xu, W., Tam, S., Zeng, X., Liu, Z., and Yin, Z. Deepthinkvla: Enhancing reasoning capability of vision-language-action models. *arXiv preprint arXiv:2511.15669*, 2025.
- Zawalski, M., Chen, W., Pertsch, K., Mees, O., Finn, C., and Levine, S. Robotic control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024.
- Zhao, T. Z., Kumar, V., Levine, S., and Finn, C. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- Zhou, X., Tie, G., Zhang, G., Wang, H., Zhou, P., and Sun, L. Badvla: Towards backdoor attacks on vision-language-action models via objective-decoupled optimization. *arXiv preprint arXiv:2505.16640*, 2025a.
- Zhou, Z., Xiao, Z., Xu, H., Sun, J., Wang, D., and Zhang, J. Goal-oriented backdoor attack against vision-language-action models via physical objects. *arXiv preprint arXiv:2510.09269*, 2025b.
- Zhu, W., Ji, X., Cheng, Y., Zhang, S., and Xu, W. {TPatch}: A triggered physical adversarial patch. In *32nd USENIX Security Symposium (USENIX Security 23)*, pp. 661–678, 2023.
- Zitkovich, B., Yu, T., Xu, S., Xu, P., Xiao, T., Xia, F., Wu, J., Wohlhart, P., Welker, S., Wahid, A., et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pp. 2165–2183. PMLR, 2023.

A. Additional Related Work

A.1. Backdoor Attacks against VLAs

Backdoor attacks constitute another important security threat to VLAs, as compromised models may behave normally on clean inputs but execute malicious behaviors once specific triggers are activated. Existing studies can be broadly categorized into untargeted (Zhou et al., 2025a; Guo et al., 2026) and targeted backdoor attacks (Zhou et al., 2025b; Li et al., 2025; Xu et al., 2025b). Untargeted backdoor attacks mainly aim to disrupt task execution and degrade policy performance under triggered conditions. For example, BadVLA (Zhou et al., 2025a) implants hidden trigger-policy associations into VLA models and induces conditional control deviations while largely preserving performance on clean inputs. In contrast, targeted backdoor attacks aim to induce attacker-specified actions or goals under triggered conditions. For example, GoBA (Zhou et al., 2025b) employs physical objects as triggers to induce goal-oriented malicious behaviors, while BackdoorVLA (Li et al., 2025) steers the model toward attacker-specified long-horizon action sequences.

In terms of attack effects, the closest line of work to ours is targeted backdoor attacks, which aim to induce attacker-specified behaviors once the trigger is activated. Nevertheless, our method is fundamentally different: rather than injecting a backdoor through model poisoning or fine-tuning, we formulate the problem as an adversarial attack that only requires optimizing an adversarial patch, without modifying the victim VLA model. This leads to a more realistic threat model in practice, as it removes the assumption of access to the victim VLA’s training pipeline. In addition, our method is more flexible, since different target behaviors can be achieved simply by changing the optimization objective and re-optimizing the patch.

A.2. CoT Attack

With the advent of reasoning LLMs, CoT has emerged as an important attack surface, as adversaries can manipulate the reasoning process to degrade model performance or circumvent safety alignment. Existing works can be broadly categorized into reasoning-corruption attacks (Xu et al., 2024; Lu et al., 2025a) and reasoning-hijacking jailbreak attacks (Yao et al., 2025; Kuo et al., 2025). The former disrupt intermediate reasoning and reduce the faithfulness or robustness of final outputs. For example, Preemptive Answer Attacks (Xu et al., 2024) show that exposing models to candidate answers before reasoning can substantially impair CoT reasoning performance, while FicDetail (Lu et al., 2025a) suggests that CoT may still be exploited to elicit harmful outputs despite its apparent safety benefits. The latter directly manipulate the reasoning trajectory or safety reasoning mechanism to induce unsafe behaviors. For instance, Mousetrap (Yao et al., 2025) injects iterative chaos into the reasoning chain to mislead large reasoning models, while H-CoT (Kuo et al., 2025) hijacks the model’s CoT-based safety reasoning process to improve jailbreak effectiveness.

Compared with reasoning LLMs, however, CoT in VLAs is typically more structured and grounded in embodied task-related elements, such as target object bounding boxes and predicted trajectories. Moreover, ECOT-lite (Chen et al., 2025a) suggests that CoT in VLAs primarily improves representation learning, rather than merely introducing intermediate reasoning steps. Therefore, existing CoT attacks developed for LLMs, which mainly rely on manipulating prompt context to steer long-horizon reasoning, may not directly transfer to the reasoning VLA.

B. Preliminary Analysis Details

B.1. Details of Active Interventions

To investigate the causal role of intermediate CoT in action generation, we perform intervention-based analysis during evaluation in simulation. We consider two intervention settings.

Instruction Masking. We first use the VLA model to generate CoT tokens conditioned on the language instruction and image observations. We then reuse the generated CoT tokens for action generation together with the original inputs, while masking the instruction tokens. This design removes the direct contribution of the instruction tokens while preserving the original positional encoding structure.

Cross-sample Shuffling. We first use the VLA model to generate CoT tokens conditioned on instruction A and image observations. We then pair these CoT tokens with instruction B and image for action generation. In this way, the semantic content of the CoT is intentionally mismatched with the instruction-observation pair used for action prediction.

In both settings, the generated action is executed in the simulation environment to obtain the next observation, and the same intervention procedure is applied at subsequent time steps in an autoregressive manner. We evaluate how these interventions

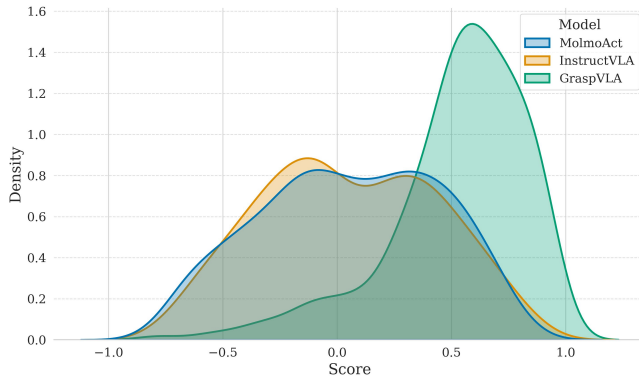


Figure 8. Visualization of hijacking score distribution across different VLAs under the shuffling setting.

affect the VLA model’s task performance, thereby assessing the causal contribution of intermediate CoT to action generation.

B.2. Further Results Analysis

As defined in Section 6.1, the score metric quantifies the relative similarity of a query trajectory to two reference trajectories. The score ranges from -1 to 1 . A value closer to 1 indicates that the query trajectory is more similar to the target trajectory T_A , while a value closer to -1 indicates greater similarity to the original trajectory T_B . A value near zero means that the query trajectory is similarly close to both reference trajectories and thus does not show a clear behavioral preference. Therefore, the sign of the score indicates the direction of alignment, and its magnitude reflects the strength of that alignment. In our experiments, we use this metric to analyze whether, under the shuffling or attack setting, the behavior of the VLA model is more strongly governed by the instruction or by the CoT. In our implementation, a higher score indicates that the generated behavior is more strongly influenced by the CoT, whereas a lower score indicates that the instruction plays a more dominant role.

We further visualize the distribution of hijack scores across different VLA models in Fig. 8. We observe that InstructVLA and MolmoAct exhibit relatively symmetric score distributions centered around zero across different tasks and layouts. This pattern suggests that, under conflicting instruction and CoT signals, the effects of the instruction and the CoT on action generation are comparatively balanced. In contrast, GraspVLA exhibits a pronounced peak in the positive-score region, indicating that its behavior is more consistently aligned with the CoT. This observation suggests that GraspVLA captures a substantially stronger causal relationship between intermediate CoT and action generation.

C. Experimental Details

C.1. VLA Model Details

To comprehensively evaluate the generalization and effectiveness of our proposed attack, we select three representative reasoning VLA models with diverse architectures, CoT mechanisms, and action decoding paradigms. Specifically, we focus on how these models formulate their intermediate reasoning steps and how they decode the final physical actions (*i.e.*, discrete tokenization vs. continuous regression). The detailed configurations of the evaluated models are as follows:

- **MolmoAct** (Lee et al., 2025): MolmoAct is an *integrated* reasoning VLA model. Built upon the Molmo, a pre-trained Vision-Language Model (VLM), it autoregressively generates depth perception tokens and visual reasoning trace tokens. Subsequent to this spatial CoT reasoning phase, it employs a *discrete* action decoding mechanism, quantizing actions into discrete bins and predicting them as low-level tokens. In implementing the proposed attack, we exclusively hijack the generation of visual reasoning trace. This targeted strategy is motivated by the fact that depth perception tokens capture task-agnostic features, whereas visual reasoning trace tokens explicitly encode task-specific intentions.
- **GraspVLA** (Deng et al., 2025): GraspVLA is an *integrated* reasoning VLA model designed for open-vocabulary grasping in continuous action spaces. Its architecture comprises an autoregressive vision-language backbone with a flow-matching action expert. Specifically, given a language instruction and images from both the front view and side

view, the model first predicts 2D target bounding boxes and then generates a grasp pose via the VLM’s autoregressive decoding. Subsequently, the original input and CoT tokens condition the generation of the action chunk through a cross-attention mechanism. This design establishes GraspVLA as a representative integrated reasoning model characterized by object-centric CoT and continuous action regression. For the proposed attack, we isolate and hijack only the 2D bounding box tokens, motivated by the practical ease of acquiring bounding box information.

- **InstructVLA** (Yang et al., 2025): InstructVLA is a *hierarchical* reasoning VLA model, but its reasoning and action generation are realized within a unified model rather than through two separate models serving as the planner and the policy, respectively. Specifically, given the instruction and image observation, the model first performs asynchronous auto-regressive reasoning to generate a textual subtask decomposition, which serves as its CoT. It then reuses the generated CoT together with the original multimodal inputs to predict latent actions through the same model, enabled by a mixture-of-experts (MoE) adaptation that allows the model to alternate between reasoning and latent action prediction. Finally, a flow-matching action expert decodes the latent actions into *continuous* action. This design establishes InstructVLA as a representative hierarchical reasoning model characterized by text-based subtask CoT and continuous action generation. For the proposed attack, we focus on hijacking the generated subtask reasoning tokens, since they explicitly encode task-level planning intentions and directly condition the subsequent latent action prediction.

C.2. Simulation Environments and Datasets

We conduct evaluation on two mainstream VLA benchmarks, **SimplerEnv** and **LIBERO**. Because the released official checkpoints differ across VLA models, InstructVLA and MolmoAct are evaluated in SimplerEnv, whereas GraspVLA is evaluated in LIBERO.

These benchmarks provide object assets and task-completion verification, which makes them suitable as standardized evaluation platforms. However, their original tasks generally contain only a limited number of objects. Since our focus is on evaluating behavior hijacking, such task settings are not sufficiently expressive. Therefore, we further curate additional tasks based on these simulation platforms and their asset libraries.

For **SimplerEnv**, we construct 10 object-centric manipulation tasks covering three motion primitives: *Pick*, *Place*, and *Move*. Each task is specified by an object pair and a corresponding instruction, enabling us to assess whether the patch can redirect the VLA from the intended behavior to an attacker-specified behavior. To control for the possibility that the attack simply induces the VLA to grasp instruction-irrelevant objects, we include at least one additional distractor object in each task. All 10 tasks are used in the preliminary analysis, while 5 representative tasks are selected for the main adversarial evaluation, as summarized in Table 7.

For **LIBERO**, we construct 10 tasks based on the implementation of the GraspVLA playground¹. As GraspVLA is designed specifically for grasping, we restrict our study to the *pick* motion primitive. Each task consists of a pair of target objects, and the corresponding instruction follows the template “pick up” + object name. In implementation, some tasks are derived directly from the original LIBERO scene layouts with modified success criteria, while other playground tasks are generated by sampling objects from the asset library and placing them randomly in the scene. Details of all tasks are provided in Table 8.

We assume a realistic threat model in which the attacker has partial prior knowledge of the task-relevant objects, but cannot determine or control their exact spatial layouts at attack deployment time. As such, the adversarial patch is expected to generalize across different object layouts rather than overfit to a fixed tabletop configuration. To simulate this setting, we randomize the object layouts for each task, using 25 layouts for optimization and 10 unseen layouts for generalization to unseen object arrangements, as illustrated in Fig. 9.

Dataset Collection. Since each curated task is defined by an object pair, we collect clean offline rollouts by executing the victim VLA under both instructions associated with that pair. Following Section 5.1, each rollout is represented as $\tau = (O, R, a)$. During optimization, the instruction specifies one object in the pair, and the target CoT/action pair (R^*, a^*) is taken from the clean rollout corresponding to the other object. In this way, the target supervision is derived from paired clean model rollouts. We further restrict target construction to the training layouts only, and reserve unseen layouts exclusively for evaluation.

¹GraspVLA playground simulation repository

Table 7. Overview of curated SimplerEnv tasks

Task	Object Pairs	Instruction	Selected for Adv.	Motion Primitive
1	green cube & yellow cube	pick green cube pick yellow cube	✓	Pick
2	coke can & pepsi can	put coke can on plate put pepsi can on plate		Place
3	apple & orange	pick apple pick orange	✓ ✓	Pick
4	apple & coke can	put apple on plate put coke can on plate		Place
5	carrot & eggplant	pick carrot pick eggplant		Pick
6	sponge & apple	pick sponge pick apple		Pick
7	coke can & orange	pick orange pick coke can	✓	Pick
8	table cloth & plate	put apple on table cloth put apple on plate	✓	Place
9	pepsi can & apple	move coke can near pepsi can move coke can near apple		Move
10	coke can - plate & orange - table cloth	put coke can on plate put orange on table cloth		Place

C.3. Instruction Analysis Details

In this section, we provide the detailed instruction configurations used in our extended evaluations. Specifically, the instruction variations crafted are listed in Table 9.

C.4. Occlusion Analysis

We further evaluate the performance of the optimized adversarial patch in simulation under occlusion settings. Specifically, we consider two occlusion strategies, center occlusion and edge occlusion, and also account for different occlusion area ratios. As shown in Table 10, the adversarial patch is more sensitive to center occlusion, and its performance degrades as the occlusion area ratio increases, which is consistent with the broader adversarial-patch literature. A possible way to improve the robustness of adversarial patches against occlusion is to actively incorporate occlusion as a form of data augmentation during optimization.

C.5. Universal Patch on Different Instructions

We further investigate the feasibility of optimizing a universal adversarial patch. Our goal is to examine whether a single patch can hijack the VLA model to execute a consistent, predefined behavior across a diverse set of instructions. To this end, we curate a fixed pool of diverse instructions for the optimization process. In each iteration, an instruction is randomly sampled to optimize the patch, thereby inducing a target behavior from the VLA model across varied linguistic contexts. Given the multimodal nature of the VLA’s decision-making process, we design two experimental settings categorized by the degree of relevance between the instruction set and the environmental objects. For example, in a scene containing an orange and an apple, a scene-relevant instruction might be “grasp the apple”. Conversely, a scene-irrelevant instruction involves semantically disjoint tasks or non-existent entities, such as “pick up the box” in a scene where no box is present. The diverse instructions utilized for evaluating the “Universal Patch” (categorized into Scene-Irrelevant and Scene-Relevant settings) are detailed in Table 11.

As illustrated in Table 12, finding a universal patch presents significant challenges when conditioned on scene-relevant

Table 8. Overview of curated LIBERO tasks

Task ID	Object Pairs	Instruction	Selected for Adv.
libero_10_0	tomato sauce can & milk	pick up tomato sauce can pick up milk	✓
libero_10_3	akita black bowl & wine bottle	pick up akita black bowl pick up wine bottle	
libero_10_4	black book & white yellow mug	pick up black book pick up white yellow mug	
libero_goal_1	akita black bowl & wine bottle	pick up akita black bowl pick up wine bottle	✓
libero_object_2	salad dressing bottle & milk	pick up salad dressing bottle pick up milk	✓
libero_object_3	alphabet soup can & chocolate pudding	pick up alphabet soup can pick up chocolate pudding	
playground_0	lemon & school bus	pick up lemon pick up school bus	
playground_1	pineapple & icecream	pick up pineapple pick up icecream	✓
playground_2	boxing glove & apple	pick up boxing glove pick up apple	
playground_3	watermelon & carrot	pick up watermelon pick up carrot	✓

Table 9. Detailed instructions used for instruction analysis.

Setting	Instructions
Extra-Context	carefully pick the orange; please pick orange; go ahead and pick orange; pick the orange on the table; pick orange in front of you
Paraphrasing	could you pick up the orange?; grasp the orange; grab the orange; pick the orange; pick up the orange

instructions. In contrast, performance metrics exhibit a substantial increase under scene-irrelevant conditions. These results suggest a synergistic influence of visual features and instructions on the efficacy of adversarial attacks.

C.6. Exploration of CLIP-based Content Loss

In addition to constraining the adversarial patch with a reference image in Section 5.3, we further explore a CLIP-based content loss to guide the semantic texture of the adversarial patch for improved visual stealthiness. As illustrated in Fig. 10, we feed a natural language prompt p (e.g., “a realistic flower”) into the CLIP text encoder and align the adversarial patch δ with the prompt in the CLIP embedding space. The CLIP-based content loss is formulated as:

$$\mathcal{L}_{\text{clip}} = 1 - \mathbf{z}_{\delta}^{\top} \mathbf{z}_p, \quad (13)$$

where \mathbf{z}_{δ} and \mathbf{z}_p denote the normalized CLIP image embedding of the adversarial patch δ and the normalized CLIP text embedding of the text prompt p , respectively. In practice, we compute $\mathcal{L}_{\text{clip}}$ over both a global view and randomly cropped local views of the patch, encouraging the optimized patch to preserve global semantic consistency while maintaining locally coherent textures.

As shown in Fig. 10, adversarial patches optimized with the CLIP-based content loss and TV regularization exhibit recognizable semantic patterns that are consistent with the given text prompts, such as flower-like or tree-like textures, while maintaining relatively smooth visual appearances.

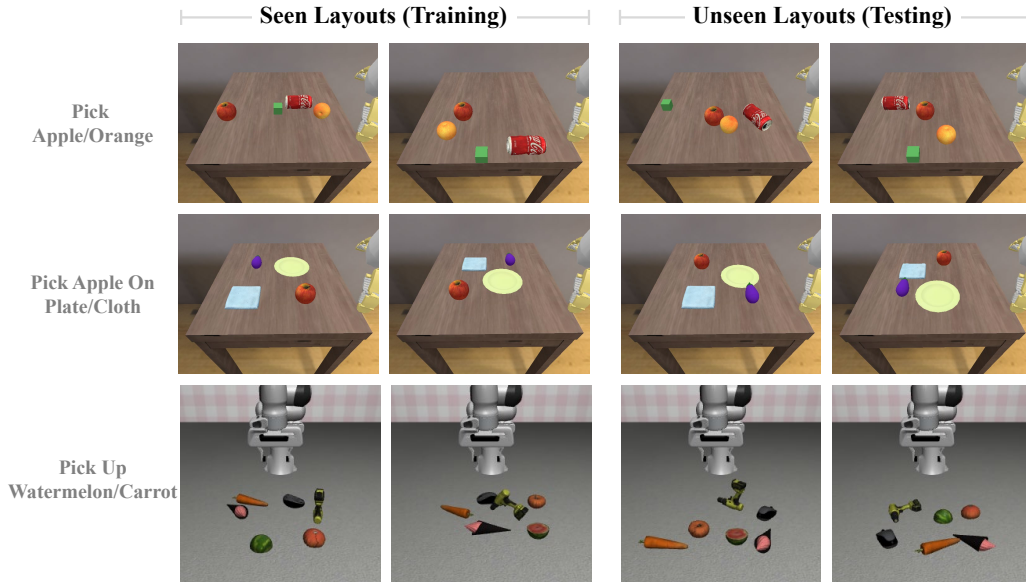


Figure 9. Representative layout variations across tasks.

Table 10. Performance under different occlusion settings.

Type	Ratio	MolmoAct		InstructVLA		GraspVLA	
		ASR (%) \uparrow	Score \uparrow	ASR (%) \uparrow	Score \uparrow	ASR (%) \uparrow	Score \uparrow
Edge	10%	44.4	0.221	62.9	0.505	8.6	-0.251
	25%	0.0	-0.363	17.1	-0.188	0.0	-0.406
	50%	0.0	-0.369	0.0	-0.358	0.0	-0.405
Center	10%	2.8	-0.320	20.0	-0.094	0.0	-0.410
	25%	0.0	-0.379	8.6	-0.162	0.0	-0.427
	50%	0.0	-0.348	0.0	-0.260	0.0	-0.408

C.7. Further Results Analysis

To further investigate the effectiveness of **TRAP** in the representation space, we visualize the latent features of the VLM in InstructVLA using t-SNE, as shown in Fig. 11. The representations under $\text{TRAP}_{\text{CoT-only}}$ attack are clearly separated from benign features, showing that the attack substantially perturbs the original latent space. Such a shift is consistent with the observed action-level failure mode, where the model repeatedly produces a single action pattern. In contrast, the representations induced by **TRAP** remain much closer to the benign feature manifold, while still maintaining strong behavior-hijacking capability.

D. Real-World Evaluation Details

D.1. Real-World Setup

As illustrated in Fig. 12, we use a Franka Panda robotic arm equipped with a parallel gripper and two Intel RealSense D415 cameras mounted at the front and side views, respectively, following the official GraspVLA real-world setup². In the simulation-aligned printed-patch setting, the optimized adversarial patch is printed at 20 cm \times 20 cm and placed on the tabletop as a standalone planar patch. In the tablecloth-style natural deployment setting, the patch is enlarged to 57 cm \times 43 cm and used as a tablecloth, with objects placed on top to introduce realistic clutter and partial occlusion.

²GraspVLA real-world controller repository

Table 11. Detailed instruction settings for universal patch analysis across different models.

Model	Setting	Instructions
InstructVLA	Scene-Irrelevant	move the yellow object; pick up the red block; place the item on the left; put the object in the container
	Scene-Relevant	grasp coke can; move orange near coke can; pick cube; pick orange
MolmoAct	Scene-Irrelevant	move the yellow object; pick up the red block; place the item on the left; put the object in the container
	Scene-Relevant	grasp apple; move orange near coke can; pick orange; pick pepsi can
GraspVLA	Scene-Irrelevant	pick up apple; pick up box; pick up nutcracker; pick up pineapple
	Scene-Relevant	pick up carrot; pick up drill; pick up mouse; pick up pumpkin

Table 12. **Universal patch analysis:** Performance under Scene-Relevant and Scene-Irrelevant instruction variations.

Victim Model	Scene-Rel.		Scene-Irrel.	
	ASR (%)	Sim.	ASR (%)	Sim.
MolmoAct	0.7	0.10	35.7	0.25
InstructVLA	25.0	0.23	65.0	0.45
GraspVLA	6.0	0.08	21.0	0.14

D.2. Color Calibration for Physical Patch

Due to distortions introduced by the printing process, camera capture, and environmental illumination, a patch optimized solely in the digital domain generally fails to retain its adversarial effectiveness once physically printed. As illustrated in Fig. 13, directly applying the adversarial patch to the image during optimization results in a pronounced sim-to-real gap, as evidenced by the comparisons in Figures 13a and 13b, thereby substantially degrading the effectiveness of the physical attack.

Therefore, to improve the visual consistency between digital-domain attack optimization and physical-domain deployment, we introduce a lightweight MLP-based color calibration module before applying the homography transformation. Specifically, we construct a reference color board, denoted as \mathcal{B}_d , and capture its printed counterpart, \mathcal{B}_f , under real-world conditions. The resulting dataset $\{\mathcal{B}_d, \mathcal{B}_f\}$ is used to train the MLP to learn the mapping between these spaces. The architecture of this module is summarized in Table 13. As shown in Fig. 13, incorporating color calibration makes the synthesized image more closely resemble its real-world counterpart, thereby facilitating the physical realization of the attack.

D.3. Impact of TV Loss Weight

We study the effect of the TV loss weight λ_3 on the stealthiness of the adversarial patches, with a particular focus on high-frequency artifacts under PGD optimization. As reported in Table 14, increasing λ_3 consistently suppresses high-frequency components in the generated patches. Specifically, the high-frequency ratio is reduced by 98.7% as λ_3 increases from 0.5 to 100. These results show that stronger TV regularization promotes smoother patches and reduces high-frequency artifacts, improving stealthiness. Moreover, we verify that the generated patches retain their attack effectiveness while maintaining high stealthiness.

D.4. Additional Analysis for Occlusion-free Deployment

Our real-world evaluation focuses on a hazardous redirection scenario, where the benign user instruction is “pick up carrot”, but the adversary aims to hijack the VLA to “pick up knife”. Across 15 independent trials, the attack successfully subverted the model’s internal reasoning for at least one step in 13/15 (86.7%) of the cases. Furthermore, in 5/15 (33.3%) of the trials, the adversary achieved full behavior hijacking, successfully completing the malicious objective from start to finish. Fig. 14 presents the synchronized multi-view observations (Front View and Side View) of this behavior hijacking process.

Failure Cases Study. We further examine the failure cases and identify two primary causes: out-of-distribution scenes and the attraction effect induced by the adversarial patch, as shown in Fig. 15.



Figure 10. Illustration of the CLIP-based objective for adversarial patch optimization. (a) The CLIP loss aligns the image embedding of an adversarial patch with the embedding of a target text prompt. (b,c) Examples of optimized adversarial patches for the prompts “a realistic flower” and “a photo of tree”, respectively.

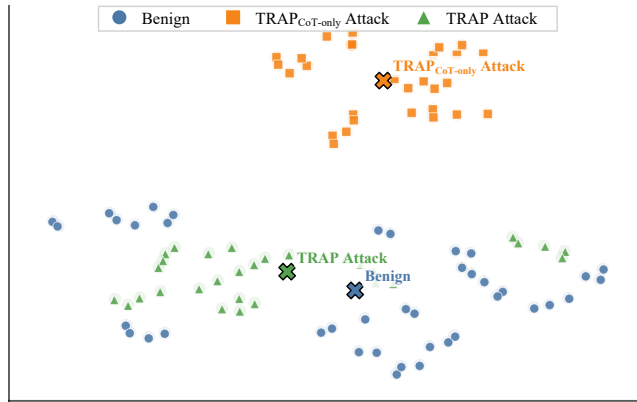


Figure 11. t-SNE visualization of latent features.

First, although the patch demonstrates a certain degree of scene-level generalization, its effectiveness degrades when the observed images differ substantially from those seen during optimization. This limitation becomes particularly pronounced under large robot-arm motions, which introduce significant visual changes into the observation. Moreover, for VLA models that use flow matching as the action expert, action generation is inherently multimodal due to the randomness of the initial noise. This additional stochasticity further increases the difficulty of reliably transferring the optimized patch across scenes and observations.

Second, we observe that the patch may exert a strong attraction effect on CoT prediction. In several cases, the bounding box predicted by GraspVLA for the target object incorrectly covers the patch itself rather than the intended object. One possible explanation is that, during optimization, the patch gradually acquires texture cues similar to those of the target object (e.g., the green knife), thereby confusing the VLA.

D.5. Additional Analysis of Object-occluded Deployment

We further evaluate **TRAP** under a more natural tablecloth-style deployment, where the optimized adversarial patch is enlarged and used as a tablecloth, with objects placed directly on top of it. Compared with the simulation-aligned printed-patch setting, this deployment is substantially more challenging: the patch is no longer an isolated planar pattern, but is partially occluded by task-relevant objects and distractors, while also being expected to remain visually inconspicuous in the scene.

The results show that **TRAP** remains effective under this more realistic physical setting. In particular, the tablecloth-style patch can still hijack the VLA behavior toward the adversarial target while preserving a semantically meaningful appearance. We further explore the robustness of the patch under several deployment variations, including different object layouts (Fig. 16a), additional distractor objects (Fig. 16b), and different lighting conditions (Fig. 16c). Across these variations, the patch consistently induces the target behavior, suggesting that the attack does not rely on a fixed tabletop configuration.

In addition to the “carrot-to-knife” hazardous-redirection task studied in the main experiment, we further evaluate another harmful object-related scenario. In this case, the benign user instruction is “pick up mouse”, while the adversary aims to

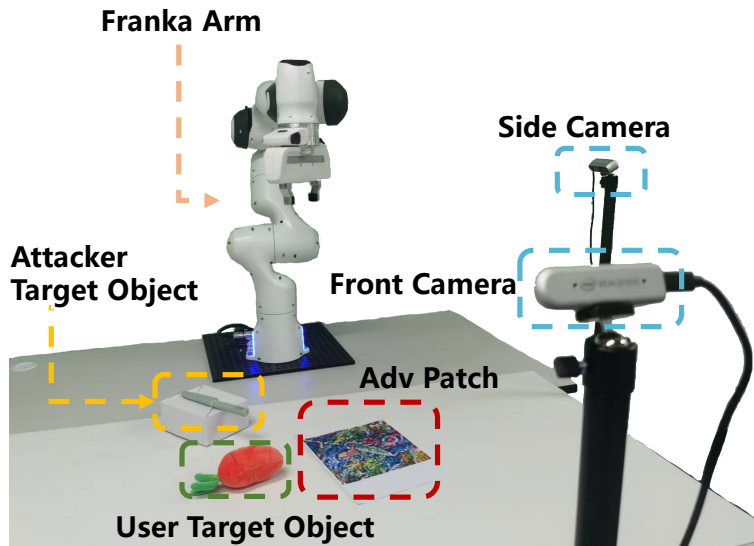


Figure 12. Real-world experiment setup.

Table 13. Architecture of the color calibration network.

Layer	Type	Input Dim.	Output Dim.
1	Linear	3	16
2	LeakyReLU	16	16
3	Linear	16	16
4	LeakyReLU	16	16
5	Linear	16	3

hijack the robot toward “pick up scissors”. As shown in Fig. 17, **TRAP** can also induce this harmful target behavior, further demonstrating its effectiveness in tablecloth-style physical deployments.

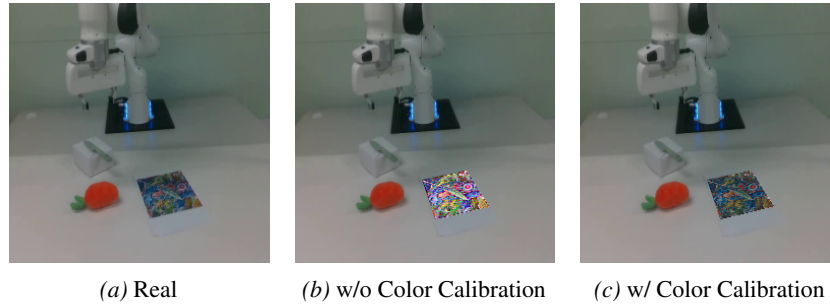


Figure 13. **Illustration of the effect of color calibration.** We compare the real image, the output without color calibration, and the output with color calibration. Color calibration reduces the visual gap between the generated output and the real observation, which is critical for real-world effectiveness.

Table 14. **Results of TV regularization on high-frequency artifacts.** Relative drop is computed with respect to $\lambda_3 = 0.5$.

TV Loss weight λ_3	0.5	1	10	100
High-frequency ratio \downarrow	0.283	0.253	0.016	0.004
Relative drop(%) \uparrow	0.0	10.6	94.3	98.7

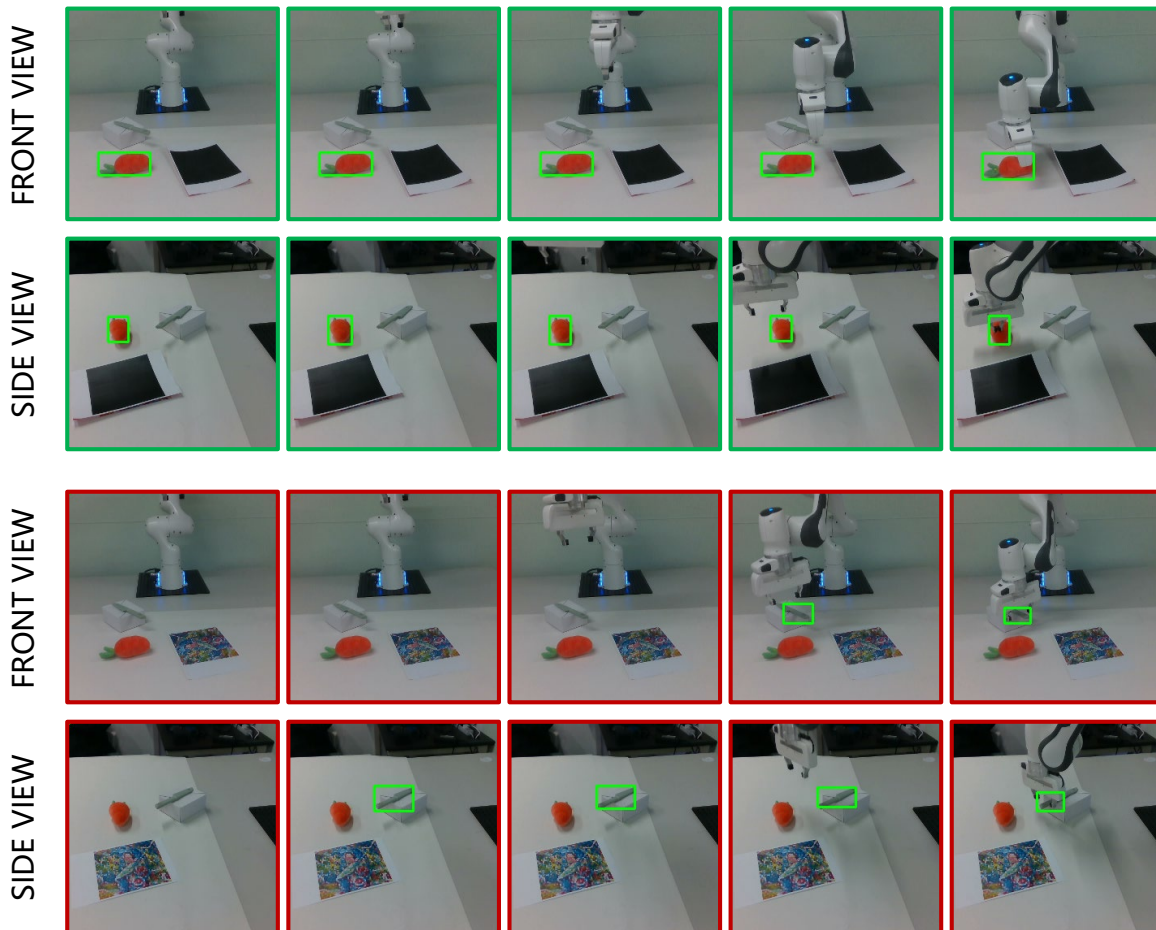


Figure 14. Detailed occlusion-free deployment result.

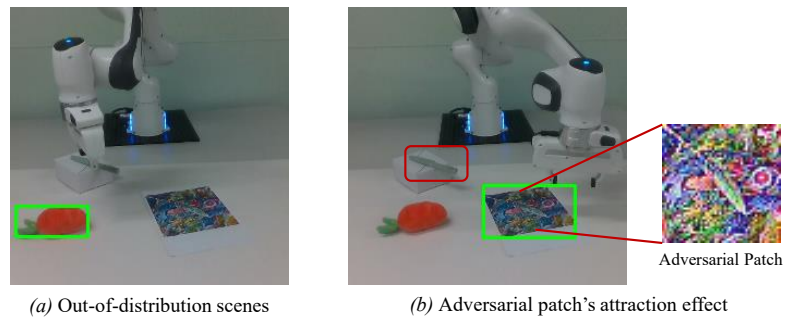
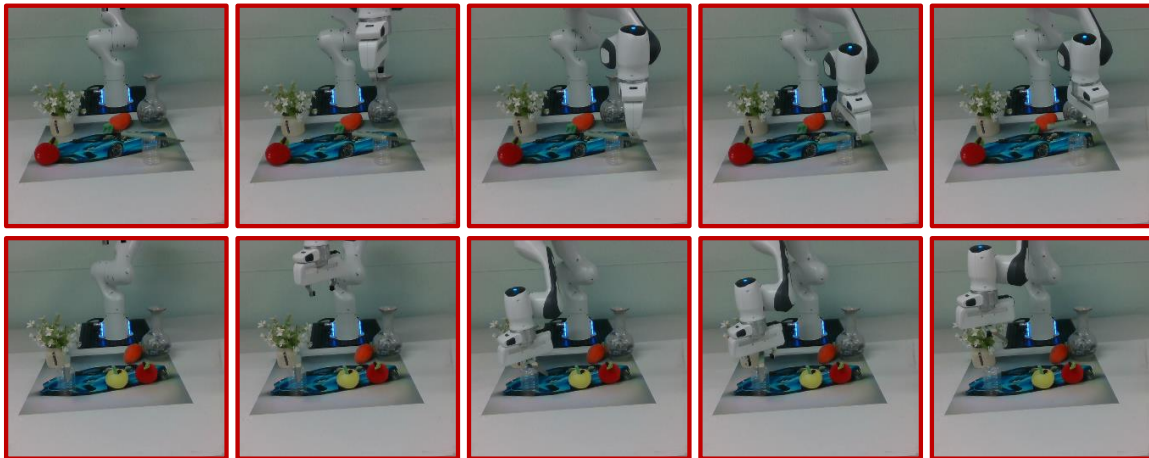
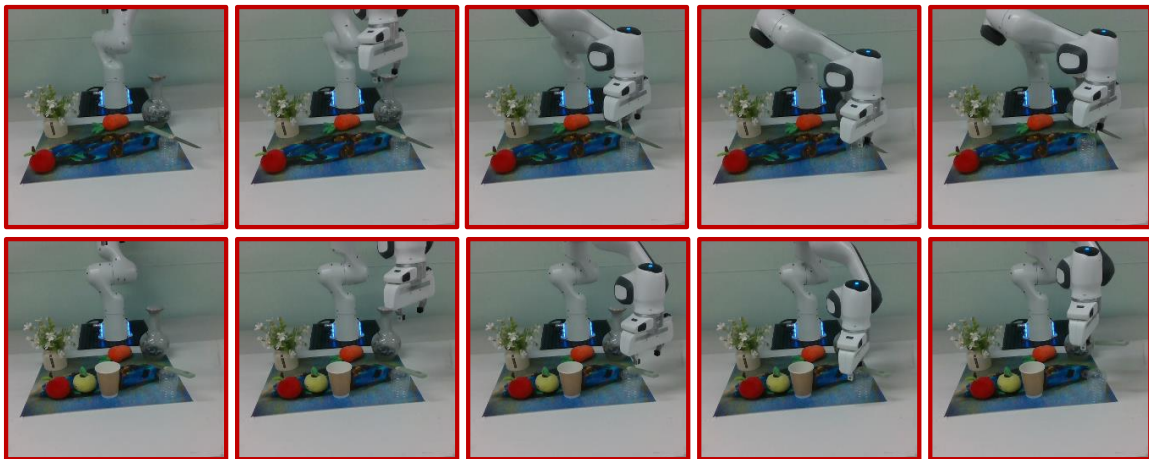


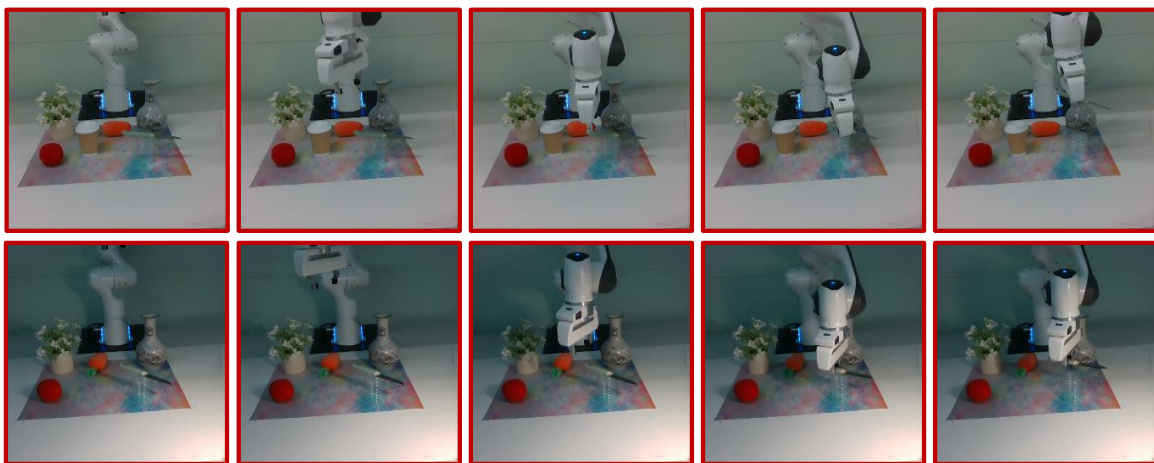
Figure 15. **Illustration of representative failure cases.** (a) Out-of-distribution observations reduce the transferability of the optimized adversarial patch, especially when large robot-arm motions introduce substantial visual changes. (b) The adversarial patch may induce an attraction effect on CoT prediction, causing GraspVLA to localize the patch itself rather than the intended target object.



(a) Evaluation under different layouts.



(b) Evaluation under add/remove objects



(c) Evaluation under different lighting conditions.

Figure 16. Illustration of evaluation under different deployment variants.



Figure 17. **Illustration of another hazardous task.** The benign user instruction is “pick up mouse”, while the adversary aims to hijack the robot toward “pick up scissors”.