

# When Robots Obey the Patch: Universal Transferable Patch Attacks on Vision-Language-Action Models

Hui Lu<sup>1</sup> Yi Yu<sup>1\*</sup> Yiming Yang<sup>1</sup> Chenyu Yi<sup>1</sup> Qixin Zhang<sup>1</sup> Bingquan Shen<sup>2</sup>  
Alex C. Kot<sup>1</sup> Xudong Jiang<sup>1</sup>

<sup>1</sup>Nanyang Technological University <sup>2</sup>DSO National Laboratories

## Abstract

*Vision-Language-Action (VLA) models are vulnerable to adversarial attacks, yet universal and transferable attacks remain underexplored, as most existing patches overfit to a single model and fail in black-box settings. To address this gap, we present a systematic study of **universal, transferable adversarial patches** against VLA-driven robots under unknown architectures, finetuned variants, and sim-to-real shifts. We introduce **UPA-RFAS** (Universal Patch Attack via Robust Feature, Attention, and Semantics), a unified framework that learns a single physical patch in a shared feature space while promoting cross-model transfer. UPA-RFAS combines (i) a feature-space objective with an  $\ell_1$  deviation prior and repulsive InfoNCE loss to induce transferable representation shifts, (ii) a robustness-augmented two-phase min-max procedure where an inner loop learns invisible sample-wise perturbations and an outer loop optimizes the universal patch against this hardened neighborhood, and (iii) two VLA-specific losses: Patch Attention Dominance to hijack text→vision attention and Patch Semantic Misalignment to induce image-text mismatch without labels. Experiments across diverse VLA models, manipulation suites, and physical executions show that UPA-RFAS consistently transfers across models, tasks, and viewpoints, exposing a practical patch-based attack surface and establishing a strong baseline for future defenses.*

## 1. Introduction

Vision-Language-Action (VLA) models have made significant strides, facilitating open-world manipulation [5, 6], language-conditioned planning [23], and cross-embodiment transfer [7, 67]. By coupling a visual encoder with language grounding and an action head, modern VLA models are capable of parsing free-form instructions and executing multi-step skills in both simulation and the physical world [31, 47]. Despite their potential, such multi-modal pipelines remain vulnerable to structured visual perturbations, which

can mislead perception, disrupt cross-modal alignment, and cascade into unsafe actions [8, 14, 29, 60]. This issue is particularly severe in robotics, as attacks that merely flip a class in perception can translate into performance drops, collisions, or violations of task constraints on real-world systems [40, 49, 63]. Motivated by that, we conduct a systematic study of universal and transferable adversarial patches for VLA-driven robots, where black-box conditions, varying camera poses, and domain shifts from simulation to the real world are the norm in practical robotic deployments.

Though vulnerabilities in VLAs have received growing attention [15, 40, 49, 57, 63], universal and transferable attacks remain largely under-explored. Reported patches often co-adapt to a specific model, datasets, or prompt template, and their success degrades sharply on unseen architectures or finetuned variants [24], precisely the black-box regimes that matter for safety assessment. As a result, current evaluations can overestimate security when the attacker lacks white-box access, and underestimate the risks of patch-based threats that exploit cross-modal bottlenecks [21]. Bridging this gap requires attacks that generalize across families of VLAs (e.g., OpenVLA [23], lightweight OFT variants [24], and flow-based policies such as  $\pi_o$  [6]).

We bridge the surrogate and victim gap by learning a universal patch in a shared feature space, guided by two principles: enlarge surrogate-side deviations that provably persist on the target, and concentrate changes along stable directions. An  $\ell_1$  deviation term drives sparse, high-salience shifts [10] that avoid surrogate-specific quirks, while a repulsive InfoNCE loss [11] pushes patched features away from their clean anchors along batch-consistent, high-CCA directions [39], strengthening black-box transfer. To further raise universality, we adopt a Robustness-augmented Universal Patch Attack (RAUP). The inner minimization loop learns a small, sample-wise invisible perturbation that reduces the feature-space objective around each input, emulating local adversarial training and hardening the surrogate. The outer maximization loop then optimizes a single physical patch against this hardened neighborhood with randomized placements and transformations, distill-

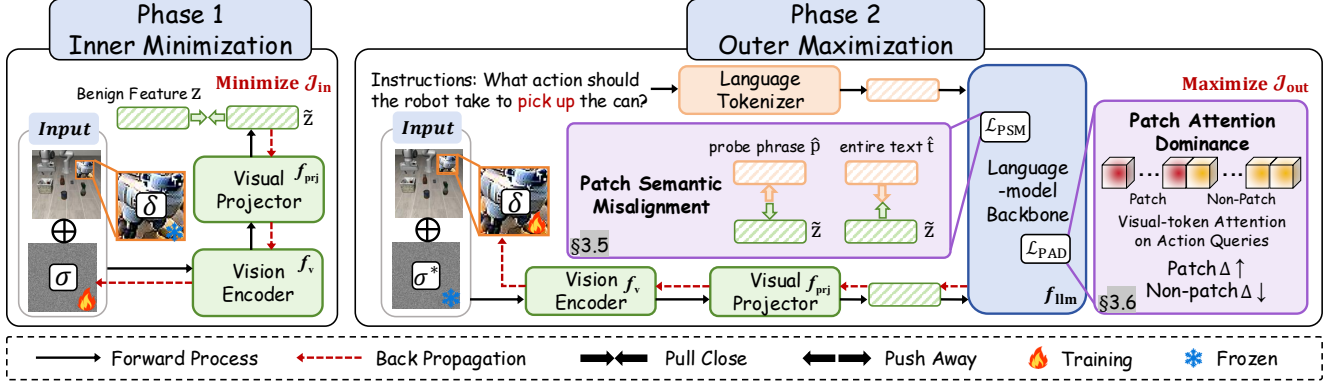


Figure 1. **Overall transferable patch attack (UPA-RFAS) for VLA robotics.** The framework operates in two coordinated stages within a shared feature-space objective. *Phase 1 – Inner minimization* learns a small, invisible, sample-wise perturbation  $\sigma$  via PGD that *minimizes* the feature objective  $\mathcal{J}_{in}$  (§ 3.3) with the patch frozen (§ 3.4). *Phase 2 – Outer maximization* freezes  $\sigma$  and optimizes a *single* physical patch  $\delta$  to *maximize*  $\mathcal{J}_{out}$  (§ 3.7), which combines an  $\ell_1$  deviation with a repulsive contrastive term and two VLA-specific objectives: **Patch Attention Dominance** (PAD) (§ 3.5) and **Patch Semantic Misalignment** (PSM) (§ 3.6). Red dashed arrows indicate back-propagation. UPA-RFAS yields a universal physical patch that transfers across models, prompts, and viewpoints.

ing the stable, cross-input directions revealed by the inner loop. For robotics, we further couple feature transfer with policy-relevant signals: (i) the Patch Attention Dominance (PAD) loss increases patch-routed text→vision attention and suppresses non-patch increments with a one-sided margin, yielding location-agnostic attention attraction; (ii) the Patch Semantic Misalignment (PSM) loss pulls the pooled patch representation toward probe-phrase anchors while repelling it from the current instruction embedding, creating a persistent image–text mismatch that perturbs instruction-conditioned policies without labels. Together, these components form **Universal Patch Attack via Robust Feature, Attention, and Semantics (UPA-RFAS)**, a universal, transferable patch framework that aligns attack feature shifts, cross-modal attention, and semantic steering.

Our contributions are summarized as follows:

- We present the first *universal, transferable* patch attack framework for VLA robotics, using a feature-space objective that combines  $\ell_1$  deviation with repulsive contrastive alignment for model-agnostic transfer.
- We propose a *robustness-augmented* universal patch attack, with invisible sample-wise perturbations as hard augmenters and a universal patch trained under heavy geometric randomization.
- We design two VLA-specific losses: *Patch Attention Dominance* and *Patch Semantic Misalignment* to hijack text→vision attention and misground instructions.
- Extensive experiments across VLA models, tasks, and sim-to-real settings show strong black-box transfer, revealing a practical patch-based threat and a transferable baseline for future defenses.

## 2. Related Work

**Vision-Language-Action (VLA) Models.** Advances in large vision–language models (LVLMs) [3, 37, 43, 59, 64]

have prompted robotic manipulation to leverage the powerful capabilities of vision–language modeling. VLA models extend LVLMs to robotic control by coupling perception, language grounding, and action generation. Autoregressive VLAs discretize actions into tokens and learn end-to-end policies from large demonstrations, yielding scalable instruction-conditioned manipulation [7, 23, 28, 38, 53, 67]. Diffusion-based VLAs generate continuous trajectories with denoisers for smooth rollouts and flexible conditioning, at the cost of higher inference latency [4–6, 27, 54]. RL-enhanced VLAs optimize robustness and adaptability beyond supervised imitation by introducing reinforcement objectives over VLA backbones [18, 34, 44]. VLA models exemplify strong vision–language alignment for compositional task understanding and end-to-end action generation, while raising new questions about robustness under instruction-conditioned deployment.

**Adversarial Attacks in Robotics.** Adversarial attacks are commonly grouped by access level: white-box methods assume full knowledge and directly use model gradients [17], whereas black-box methods operate without internals—either by querying the model for feedback [9] or by exploiting cross-model transferability of crafted examples. To strengthen transfer, optimization-driven approaches refine or stabilize gradient signals to avoid local minima arising from mismatched decision boundaries across architectures [12, 26, 30, 32, 50, 65]. Augmentation-based strategies diversify inputs to induce gradient variation and reduce overfitting to a single surrogate [13, 30, 48, 51, 56, 66]. Finally, feature-space attacks aim at intermediate representations to promote cross-model invariance and further improve transfer [16, 52, 61, 62]. Patch-based physical attacks [8, 19, 41, 55, 58] are practical for real-world deployment, remaining effective under changes in viewpoint and illumination, which makes them suitable for robotic systems.

VLA models [20, 45] couple visual and linguistic modalities to align perception with action, and visual streams are high-dimensional and can be subtly perturbed in ways that are difficult to detect [2, 46]. Accordingly, our work targets the visual modality with a universal, transferable patch attack. To our knowledge, it is the first to investigate black-box transfer vulnerabilities of VLAs in real-world settings.

### 3. Methodology

#### 3.1. Preliminary

**Adversarial Patch Attack.** We consider a robot whose decisions are based on RGB visual streams  $\mathbf{x}_t \in [0, 1]^{H \times W \times 3}$  across time step  $t$ . An adversary tampers with this stream using a *single universal* patch  $\delta \in [0, 1]^{h_p \times w_p \times 3}$ . At each time step  $t$ , an area-preserving geometric transformation  $T_t \sim \mathcal{T}$  (e.g., random position, skew, and rotation) is sampled, and the transformed patch is rendered onto the frame. Given  $\mathbf{M}_{T_t} \in \{0, 1\}^{H \times W}$  as the binary placement mask induced by  $T_t$ , and  $\mathcal{R}(\delta; T_t) \in [0, 1]^{H \times W \times 3}$  as the rendered patch, the pasting function  $\mathcal{P}$  and resulting patched frame is

$$\begin{aligned} \tilde{\mathbf{x}}_t &= \mathcal{P}(\mathbf{x}_t, \delta, T_t) \\ &= (\mathbf{1} - \mathbf{M}_{T_t}) \odot \mathbf{x}_t + \mathbf{M}_{T_t} \odot \mathcal{R}(\delta; T_t) \quad \text{s.t. } S(\delta) < \rho, \end{aligned} \quad (1)$$

where  $\odot$  is the Hadamard product,  $\mathbf{1}$  is an all-ones matrix,  $S(\cdot)$  returns the patch area (i.e.,  $h_p \times w_p$ ), and  $\rho$  is an area budget controlling the maximal visible size of the patch.

Let  $\pi$  denote a *victim* policy. Given visual inputs  $\mathbf{x}$  drawn from a task distribution  $p(\mathbf{x})$  and random patch placements  $T_t \sim \mathcal{T}$ , an adversarial patch attack aims to *learn* a single universal patch  $\delta$  that maximizes an evaluation objective  $\mathcal{J}_{\text{eval}}$  (e.g., task loss increase or action-space deviation [49]) under pasting  $\mathcal{P}$  in Eq. 1 and these randomized conditions:

$$\delta^* \in \arg \max_{S(\delta) < \rho} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \mathcal{J}_{\text{eval}}(\mathcal{P}(\mathbf{x}, \delta, T_t); \pi) \right]. \quad (2)$$

This objective captures a *single* patch that is robustly effective across time, viewpoints, and scene configurations.

**VLAs.** We follow the OpenVLA formulation [23], where a policy is decomposed into a *vision encoder*  $f_v$ , a *visual projector*  $f_{\text{prj}}$ , and a *language-model backbone*  $f_{\text{llm}}$  equipped with an *action head*  $f_{\text{act}}$ . Given an RGB observation  $\mathbf{x}$  and an instruction  $c$ , the model predicts an action vector  $\mathbf{y}$  as

$$\mathbf{y} = \text{OpenVLA}(\mathbf{x}, c) = f_{\text{act}}(f_{\text{llm}}([f_{\text{prj}}(f_v(\mathbf{x})), \text{tok}(c)])) . \quad (3)$$

The computation can be unpacked as: **(i)** the vision encoder  $f_v$  maps the image into a set of multi-granularity visual embeddings, for example by concatenating DINOv2 [37] and SigLIP [59] features, yielding  $\mathbf{E}_v \in \mathbb{R}^{N_v \times D_v}$  from  $\mathbf{x}$ ; **(ii)** the projector  $f_{\text{prj}}$  aligns these embeddings to the LLM token space, producing visual tokens  $\mathbf{Z}_v \in \mathbb{R}^{N_v' \times D_t}$ ; **(iii)** the backbone  $f_{\text{llm}}$  takes the concatenation of  $\mathbf{Z}_v$  and the tokenized command  $\text{tok}(c)$ , and fuses them into hidden states

$\mathbf{H}_\ell$ ; **(iv)** the action head  $f_{\text{act}}$  decodes  $\mathbf{H}_\ell$  into the continuous control output  $\mathbf{y} \in \mathbb{R}^{D_a}$  (e.g., a 7-DoF command).

#### 3.2. Problem Formulation

Existing VLA patch attacks [49] assume white-box access to the victim model, which limits their practicality and says little about cross-policy transfer. In our setting, the attacker instead only has gradient access to a *single* surrogate model  $\hat{\pi}$  and aims to learn one universal patch that transfers to a family of unseen target policies  $\Pi_{\text{tgt}}$ . To formalize this threat model, we separate *optimization* and *evaluation*: the patch is optimized in the surrogate feature space via a differentiable objective  $\mathcal{J}_{\text{tr}}$ , and its success is assessed by an evaluation objective  $\mathcal{J}_{\text{eval}}$  on target policies drawn from  $\Pi_{\text{tgt}}$ . Following [49], we adopt the untargeted attack setting, and summarize this transferable patch attack as follows.

**Definition 1 (Transferable adversarial patch attack via VLA feature space)** Let  $\hat{\pi}$  be a surrogate model and  $\Pi_{\text{tgt}}$  a family of target policies. Let  $f_{\hat{\pi}}(\cdot)$  extract features from  $\hat{\pi}$ . A patch  $\delta$  is a universal transferable adversarial patch in the VLA feature space, it satisfies

$$\begin{aligned} \max_{\delta_s} \quad & \mathbb{E}_{\pi \sim \Pi_{\text{tgt}}} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \mathcal{J}_{\text{eval}}(\mathcal{P}(\mathbf{x}, \delta_s, T_t); \pi) \right] \\ \text{s.t. } \delta_s \in \arg \max_{\delta} \quad & \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \mathcal{J}_{\text{tr}}(\mathcal{P}(\mathbf{x}, \delta, T_t); \hat{\pi}) \right], \end{aligned} \quad (4)$$

where  $\mathcal{J}_{\text{tr}}$  measures feature discrepancy using  $\Delta$ :

$$\mathcal{J}_{\text{tr}}(\mathcal{P}(\mathbf{x}, \delta, T_t); \hat{\pi}) = \Delta(f_{\hat{\pi}}(\mathcal{P}(\mathbf{x}, \delta, T_t)), f_{\hat{\pi}}(\mathbf{x})). \quad (5)$$

Here,  $\mathcal{P}$  is the pasting function defined in Eq. 1, and  $\mathcal{J}_{\text{tr}}$  is the transferable attack strategy. Although  $\hat{\pi}$  and  $\pi$  differ in training recipe and data, we probe whether their feature spaces admit a stable *cross-model relation* as follows:

#### Shared Representational Structure across VLA Policies.

Empirically, we observe a strong linear relationship between the surrogate and target feature spaces. Let  $\mathbf{z}_s$  and  $\mathbf{z}_t$  denote visual features from  $\hat{\pi}$  and  $\pi$  on the same inputs. We first apply Canonical Correlation Analysis (CCA) to test whether these representations lie in a *shared linear subspace*: large top Canonical Correlation indicates a near-invertible linear map aligning the two subspaces [36, 39]. In parallel, we fit a *linear regression probe* from  $\mathbf{z}_s$  to  $\mathbf{z}_t$  and use the explained variance ( $R^2$ ) to quantify how well a *single* linear map accounts for the target features, complementing CCA’s subspace view [1, 25]. In our case,  $R^2 \approx 0.654$  together with near-unity top- $k$  Canonical Correlations indicates a shared low-dimensional subspace, with some residual components not captured by one linear map. Consequently, patch updates that steer  $\hat{\pi}$ ’s features within this shared subspace tend to induce homologous displacements in  $\pi$ , supporting the transferability of patches. Motivated by these observations, we make the following Assumption 1.

### 3.3. Learning Transferable Patches with Feature-space $\ell_1$ and Contrastive Misalignment

Let  $f_{\hat{\pi}}, f_{\pi} : \mathcal{X} \rightarrow \mathbb{R}^d$  be the surrogate and target encoders with dimension  $d$ , where  $f$  consists of vision encoder  $f_v$  and visual projector  $f_{\text{prj}}$ . For any pair  $(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$ , define the surrogate-side feature deviation  $\Delta \mathbf{z}_i := f_{\hat{\pi}}(\tilde{\mathbf{x}}_i) - f_{\hat{\pi}}(\mathbf{x}_i)$  and the target-side deviation  $\Delta \mathbf{g}_i := f_{\pi}(\tilde{\mathbf{x}}_i) - f_{\pi}(\mathbf{x}_i)$ .

**Assumption 1 (Linear alignment with bounded residual)** *There exists a matrix  $A^* \in \mathbb{R}^{d \times d}$  such that*

$$f_{\pi}(\mathbf{x}) = f_{\hat{\pi}}(\mathbf{x}) A^* + e(\mathbf{x}), \quad (6)$$

where the alignment residual  $e(\mathbf{x})$  satisfies  $\|e(\tilde{\mathbf{x}}) - e(\mathbf{x})\|_2 \leq \varepsilon_E$  for all pairs  $(\mathbf{x}, \tilde{\mathbf{x}})$  considered.

Assumption 1 states that the effect of a surrogate deviation must persist on the target with strength governed by  $\sigma_{\min}(A^*)$ , the smallest singular value of the alignment map  $A^*$ . The proposition below makes this dependence explicit.

**Proposition 1 (Lower-bounded target displacement)**

*Under Assumption 1, for any  $(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$*

$$\|\Delta \mathbf{g}_i\|_2 \geq \sigma_{\min}(A^*) \|\Delta \mathbf{z}_i\|_2 - \varepsilon_E, \quad (7)$$

and, using Hölder's inequality  $\|v\|_1 \leq \sqrt{d}\|v\|_2$ ,

$$\|\Delta \mathbf{g}_i\|_1 \geq \frac{\sigma_{\min}(A^*)}{\sqrt{d}} \|\Delta \mathbf{z}_i\|_1 - \varepsilon_E. \quad (8)$$

Proof is in Appendix A. Proposition 1 links target-side deviation to the surrogate-side. Therefore, any strategy that enlarges  $\|\Delta \mathbf{z}_i\|$  (e.g., via an  $\ell_1$  objective) necessarily induces a nontrivial response on the target. Thus we can capture why we could use **L1 loss**  $\mathcal{L}_1$  to maximize feature discrepancy:

**Corollary 1 (Effect of maximizing  $\ell_1$  deviation)** *If an attack increases the surrogate-side  $\ell_1$  deviation, e.g. by maximizing  $\mathcal{L}_1 = \|\Delta \mathbf{z}_i\|_1$ , then the target-side deviation obeys the linear lower bound in Eq. 8. In particular, when the alignment is well-conditioned ( $\sigma_{\min}(A^*)$  not small) and the residual coupling  $\varepsilon_E$  is modest, increasing  $\|\Delta \mathbf{z}_i\|_1$  necessarily induces a nontrivial increase of  $\|\Delta \mathbf{g}_i\|_1$ .*

**Repulsive Contrastive Regularization.** Complementing the  $\mathcal{L}_1$  deviation term, we introduce a *repulsive* contrastive objective that explicitly pushes the patched feature  $\tilde{\mathbf{z}}_i$  away from its clean anchor  $\mathbf{z}_i$ . For each sample  $i$ , we still treat  $(\mathbf{z}_i, \tilde{\mathbf{z}}_i)$  as a distinguished pair and  $\{\tilde{\mathbf{z}}_j\}_{j \neq i}$  as a reference set, and adopt the InfoNCE [11] as a repulsion term

$$\mathcal{L}_{\text{con}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_j)/\tau)}, \quad (9)$$

where  $\text{sim}$  denotes cosine similarity and  $\tau$  is a temperature. Maximizing (minimizing)  $\mathcal{L}_{\text{con}}$  encourages the similarity  $\text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i)$  to *decrease (increase)*, effectively pushing

$\tilde{\mathbf{z}}_i$  away (pulling  $\tilde{\mathbf{z}}_i$  close) from its clean anchor and concentrating the change along directions that are consistently shared across the batch.

**Overall Feature-space Objective.** Combining both components, we obtain the objective for  $\Delta$  as given in Eq. 5:

$$\mathcal{J}_{\text{tr}} = \mathcal{L}_1 + \lambda_{\text{con}} \mathcal{L}_{\text{con}}, \quad (10)$$

where  $\mathcal{L}_1$  is the  $\ell_1$  loss term and  $\mathcal{L}_{\text{con}}$  is the repulsive contrastive objective, and  $\lambda > 0$  balances their contributions.

### 3.4. Robustness-augmented Universal Patch Attack

**Emulate Robust Surrogates without Retraining VLAs.**

Transfer-based attacks on image classifiers have shown that adversarial examples generated on *adversarially trained* or *slightly robust* source models transfer significantly better than those crafted on standard models [22, 42]. Robust training encourages the source model to rely on more “universal” features shared across architectures, so perturbations aligned with these features exhibit stronger cross-model transferability. A natural strategy would be to use an adversarially trained VLA as the surrogate. However, adversarially trained large VLA policies is practically prohibitive: it requires massive interactive data and computing, and can substantially degrade task performance. Instead, we still optimize a *single universal physical patch*, but augment it with a *sample-wise, invisible* perturbation that *emulates* adversarial (robust) training on the surrogate. This perturbation is applied globally and updated to *counteract* patch-induced feature deviations, effectively “hardening” the surrogate along the directions the patch tries to exploit. Since the universal patch is localized while the sample-wise perturbations remain invisible and input-specific, their interference is limited, and the patch can then exploit the robust feature directions revealed by this hardening step.

**Bi-level Robustness-augmented Optimization.** Formally, let  $\delta$  denote the universal patch and  $\sigma$  a sample-wise perturbation confined to the patch mask. Given a optimizing loss  $\mathcal{J}_{\text{tr}}$  on the surrogate  $\hat{\pi}$ , we consider the following robustness-augmented bi-level objective:

$$\begin{aligned} \delta^* \in \arg \max_{\mathcal{S}(\delta) < \rho} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathcal{J}_{\text{tr}} \left( \mathcal{P}(\mathbf{x} + \sigma^*(\mathbf{x}), \delta, T_t); \hat{\pi} \right) \\ \text{s.t. } \sigma^*(\mathbf{x}) \in \arg \min_{\|\sigma\|_{\infty} \leq \epsilon_{\sigma}} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathcal{J}_{\text{tr}} \left( \mathcal{P}(\mathbf{x} + \sigma, \delta, T_t); \hat{\pi} \right). \end{aligned} \quad (11)$$

The inner problem “adversarially trains” the surrogate locally by finding a small, sample-wise perturbation to *reduce* the attack loss, and the outer problem then maximizes the same loss with respect to  $\delta$  in this hardened neighborhood.

To further strengthen transferability, the outer maximization is not driven by the feature displacement alone. In the following subsections, we introduce **additional loss components** that shape *where* the model attends and *what* semantics the patch encodes, and jointly optimize them within this robustness-augmented framework.



### 3.5. Patch Attention Dominance: Cross-Modal Hijack Loss

**Action-relevant Queries as the Attack Handle.** In VLA policies, actions are largely driven by a small set of *action-relevant* text queries whose cross-modal attention to vision decides which visual regions control the policy. Our universal patch is therefore designed as a *location-agnostic attention attractor*: regardless of placement, skew, or orientation, it should **redirect the attention** of these action-relevant queries *from true semantic regions to the patch*. Concretely, we aim to *increase* the attention increments on the patch vision tokens while *reducing* increments on non-patch tokens, based on the difference between patched and clean runs under random placements.

**Patch-induced Attention Increments for Action-relevant queries.** From clean and patched runs, we collect the last  $N$  attention blocks  $\mathbf{A}$  from  $f_{\text{llm}}$ , average over heads, and slice out the text  $\rightarrow$  vision submatrix via  $\text{tv}(\cdot)$ :

$$\begin{aligned}\bar{\mathbf{A}}_c^{(l)} &= \frac{1}{H} \sum_{h=1}^H \mathbf{A}_{c,:,:,h}^{(l)}, & \mathbf{B}_c^{(l)} &= \text{tv}(\bar{\mathbf{A}}_c^{(l)}), \\ \bar{\mathbf{A}}_p^{(l)} &= \frac{1}{H} \sum_{h=1}^H \mathbf{A}_{p,:,:,h}^{(l)}, & \mathbf{B}_p^{(l)} &= \text{tv}(\bar{\mathbf{A}}_p^{(l)}),\end{aligned}\quad (12)$$

where  $l = L - N + 1, \dots, L$  indexes the last  $N$  layers. We row-normalize over vision tokens (index  $p$ ) and average across layers to obtain attention *shares*, then define the patch-induced share increment:

$$\Delta = \frac{1}{N} \sum_l \text{rn}(\mathbf{B}_p^{(l)}) - \frac{1}{N} \sum_l \text{rn}(\mathbf{B}_c^{(l)}) \in \mathbb{R}^{B \times T \times P}, \quad (13)$$

where  $\text{rn}(\cdot)$  denotes row-normalization over  $p$ . By optimizing  $\Delta$  rather than raw attention, the objective depends only on *patch-induced* changes.

**Action-relevant Queries.** To focus precisely on action-relevant queries and avoid surrogate-specific overfitting, we restrict the optimization to the top- $\rho$  text tokens (per batch) that already receive the highest clean attention:

$$\tilde{\Delta} = \Delta \odot \chi, \quad \chi = \text{TopKMask}(\mathbf{B}_c; \rho), \quad (14)$$

where  $\text{TopKMask}$  returns a binary mask over the text positions, broadcast across vision tokens. These top- $\rho$  tokens are our proxy for action-relevant queries.

**Patch vs. Non-patch Attention Increments.** To capture the effect of patch location on visual tokens, we map the pixel-level mask  $\mathbf{M}_{T_i} \in \{0, 1\}^{H \times W}$  to a token-level mask  $\mathbf{M}_z \in [0, 1]^p$  via bilinear interpolation, where  $p$  is the number of visual tokens (e.g.,  $p = g^2$  for a  $g \times g$  ViT grid), and then flatten it to length  $p$ . We then aggregate the attention increments routed from action-relevant queries into patch versus non-patch vision tokens:

$$\begin{aligned}d_{\text{patch}} &= \langle \tilde{\Delta}, \mathbf{M}_z \rangle_p, & d_{\text{non}} &= \langle \tilde{\Delta}, \mathbf{1} - \mathbf{M}_z \rangle_p, \\ \text{non\_top} &= \max_p (\tilde{\Delta} \odot (\mathbf{1} - \mathbf{M}_z)),\end{aligned}\quad (15)$$

where  $\langle \cdot, \cdot \rangle_p$  sums over the vision index  $p$ .  $d_{\text{patch}}$  ( $d_{\text{non}}$ ) measures how much extra attention the patch induces on patch (non-patch) tokens from action-relevant text queries.

**Patch Attention Dominance (PAD) Loss.** Finally, we define the attention-hijack objective to maximize by explicitly *increasing* patch-related increments and *decreasing* non-patch increments, with a margin against the strongest non-patch route:

$$\begin{aligned}\mathcal{L}_{\text{PAD}} &= \mathbb{E}[d_{\text{patch}}] - \lambda \mathbb{E}[\text{ReLU}(d_{\text{non}})] \\ &\quad - \mathbb{E}[\text{ReLU}(m - (d_{\text{patch}} - \text{non\_top}))],\end{aligned}\quad (16)$$

where  $\mathbb{E}[\cdot]$  averages over the selected (action-relevant) text tokens. The first term increases patch attention increments, the second penalizes positive increments on non-patch tokens, and the margin term enforces that the patch’s increment exceeds the strongest non-patch increment by at least  $m$ . Together, these terms induce *Patch Attention Dominance*, where action-relevant queries direct their additional attention to the patch rather than true semantic regions.

### 3.6. Patch Semantic Misalignment: Text-Similarity Attack Loss

**Semantic Steering beyond Attention.** Merely hijacking cross-modal attention does not guarantee a consistent behavioral bias across models or tasks. To further enhance transferability, we constrain the patch also in *semantic* space: we **steer** the visual representation of patch-covered tokens **toward** a set of cross-model-stable action/direction primitives (*probe phrases*), while simultaneously **pushing it away from the holistic representation** of the current instruction. The probes (e.g., “put”, “pick up”, “place”, “open”, “close”, “left”, “right”) act as architecture-agnostic anchors, and the repulsion from the instruction embedding induces a persistent, context-dependent semantic misalignment that more reliably derails the policy decoder.

**Patch Pooling and Semantic Anchors.** Let  $\mathbf{z}_j \in \mathbb{R}^D$  be visual token features and  $m_j \in \mathbf{M}_z$  the corresponding patch-token mask. We pool and  $\ell_2$ -normalize the patch feature:

$$\hat{\mathbf{v}}_{\text{patch}} = \left\| \left( \sum_{j=1}^P m_j \mathbf{z}_j \right) / \left( \sum_{j=1}^P m_j + \varepsilon \right) \right\|_2. \quad (17)$$

Let  $\{\hat{\mathbf{p}}_k\}_{k=1}^K$  be normalized *probe prototypes* (e.g., action and direction anchors), and let  $\hat{\mathbf{t}}$  denote a normalized representation of the whole current instruction (e.g., the mean of the last-layer text states from  $f_{\text{llm}}$ ).

**Patch Semantic Misalignment (PSM) Loss.** We then define the text-similarity attack loss to maximize as

$$\mathcal{L}_{\text{PSM}} = \alpha \left[ \log \sum_{k=1}^K \exp \left( \frac{\hat{\mathbf{v}}_{\text{patch}}^\top \hat{\mathbf{p}}_k}{\tau} \right) \right] - \beta \hat{\mathbf{v}}_{\text{patch}}^\top \hat{\mathbf{t}}, \quad (18)$$

with temperatures  $\tau > 0$  and weights  $\alpha, \beta > 0$ .

Eq. 17 yields a location-agnostic semantic descriptor for the patch-covered tokens. In Eq. 18, the first (LogSumExp) term *pulls*  $\hat{\mathbf{v}}_{\text{patch}}$  toward any probe prototype, avoiding dependence on a single phrase while focusing gradients on the most compatible anchors as  $\tau$  decreases. The second term *pushes* the patch feature away from the instruction embedding, inducing a persistent, context-dependent semantic mismatch, with  $\alpha, \beta$  balancing pull and push. The loss is fully differentiable w.r.t. patch parameters via  $\mathbf{z}_j$  and complements attention hijacking by steering the *attended* content toward a stable, transferable semantic direction.

### 3.7. Universal Patch Attack via Robust Feature, Attention, and Semantics (UPA-RFAS)

The overall optimization process is in Algorithm 1 where:

**Inner Minimization.** Given  $\mathbf{x}$  at time  $t$  and the current patch  $\delta$ , we initialize a global invisible perturbation  $\sigma^{(0)} = \mathbf{0}$  and update it by Projected Gradient Descent (PGD) [35]:

$$\sigma^{(i+1)} \leftarrow \Pi_{\|\cdot\|_\infty \leq \epsilon_\sigma} \left( \sigma^{(i)} - \eta_\sigma \nabla_\sigma \mathcal{J}_{\text{in}} \left( \mathcal{P}(\mathbf{x} + \sigma^{(i)}, \delta, T_t); \hat{\pi} \right) \right), \quad (19)$$

where  $\mathcal{J}_{\text{in}} = \mathcal{J}_{\text{tr}}$  is in Eq. 10,  $\Pi_{\|\cdot\|_\infty \leq \epsilon_\sigma}$  projects onto the  $\ell_\infty$  ball of radius  $\epsilon_\sigma$ ,  $\eta_\sigma$  is the step size, and  $T_t$  is sampled once from  $\mathcal{T}$  across iterations.  $\sigma^*$  is obtained after  $I$  iterations.

**Outer Maximization.** With  $\sigma^*(\mathbf{x})$  fixed, we update the universal patch  $\delta$  by AdamW [33] to maximize the objective with additional losses under randomized transformations:

$$\delta \leftarrow \text{AdamW} \left( -\mathcal{J}_{\text{out}} \left( \mathcal{P}(\mathbf{x} + \sigma^*(\mathbf{x}), \delta, T_t); \hat{\pi} \right); \eta_\delta \right), \quad (20)$$

$$\mathcal{J}_{\text{out}} = \mathcal{L}_1 + \lambda_{\text{con}} \mathcal{L}_{\text{con}} + \lambda_{\text{PAD}} \mathcal{L}_{\text{PAD}} + \lambda_{\text{PSM}} \mathcal{L}_{\text{PSM}}, \quad (21)$$

where the patch  $\delta \in [0, 1]^{h_p \times w_p \times 3}$  respects the area budget, and  $\eta_\delta$  is the learning rate. At each iteration, we sample  $T_t \sim \mathcal{T}$  and clamp  $\delta$  to the valid range  $[0, 1]$ .

## 4. Experiments

### 4.1. Experiment Setup

**Datasets.** We evaluate our attacks on BridgeData V2 [47] and LIBERO [31] using the corresponding VLA models. BridgeData V2 is a real-world corpus spanning 24 environments and 13 manipulation skills (e.g., grasping, placing, object rearrangement), comprising 60,096 trajectories. LIBERO is a simulation suite with four task families-Spatial, Object, Goal, and Long, where LIBERO-Long combines diverse objects, layouts, and extended horizons, making multi-step planning particularly challenging.

**Baseline.** We adopt RoboticAttack [49]’s 6 objectives as the baselines, including Untargeted Manipulation Attack (UMA), Untargeted Action Discrepancy Attack (UADA), and Targeted Manipulation Attack (TMA) corresponding to

---

#### Algorithm 1 UPA-RFAS

---

```

1: Input: surrogate  $f_{\hat{\pi}}$ , subset  $\mathcal{D}_s$ , universal patch  $\delta \in [0, 1]^{h_p \times w_p \times 3}$ , budget  $\epsilon_\sigma$ , inner steps  $I$ , outer steps  $K$ , step sizes  $\eta_\sigma, \eta_\delta$ , weights  $\lambda_{\text{con}}, \lambda_{\text{PAD}}, \lambda_{\text{PSM}}$ 
2: for mini-batch data  $(\mathbf{x}, c, t) \in \mathcal{D}_s$  do
3:   # Inner minimization
4:   Initialize sample-wise perturbation  $\sigma^{(1)} \leftarrow \mathbf{0}$ 
5:   Sample  $T_t \sim \mathcal{T}$ 
6:   for  $i = 1$  to  $I$  do
7:      $\mathcal{J}_{\text{in}} = \mathcal{J}_{\text{tr}} \left( \mathcal{P}(\mathbf{x} + \sigma^{(i)}, \delta, T_t); \hat{\pi} \right)$  via Eq. 1 and 10
8:      $\sigma^{(i+1)} \leftarrow \Pi_{\|\cdot\|_\infty \leq \epsilon_\sigma} \left( \sigma^{(i)} - \eta_\sigma \nabla_\sigma \mathcal{J}_{\text{in}} \right)$ 
9:   end for
10:   $\sigma^* \leftarrow \sigma^{(I)}$ 
11:  #Outer maximization
12:  for  $k = 1$  to  $K$  do
13:    Sample  $T_t \sim \mathcal{T}$ 
14:    Compute  $\mathcal{J}_{\text{out}}$  via Eq. 1, 16, 18 and 21
15:     $\delta \leftarrow \text{AdamW} \left( -\mathcal{J}_{\text{out}} \left( \mathcal{P}(\mathbf{x} + \sigma^*(\mathbf{x}), \delta, T_t); \hat{\pi} \right); \eta_\delta \right)$ 
16:     $\delta \leftarrow \text{Clip}_{[0,1]}(\delta)$ 
17:  end for
18: return  $\delta$ 
```

---

different Degree-of-freedom (DoF). For each, experiments follow the original loss definitions and evaluation protocol. We further consider both **simulated and physical victim settings**: a model trained in simulation on the LIBERO-Long suite using the *OpenVLA-7B-LIBERO-Long* variant, and a model trained on real-world BridgeData v2 data with the *OpenVLA-7B* model, respectively.

**Surrogate and Victim VLAs.** We evaluate universal, transferable patches under a strict black-box transfer protocol. Surrogate models are chosen from publicly available, widely used VLA [23] to reflect prevailing design trends. The primary surrogate models are *OpenVLA-7B* trained on physical dataset BridgeData V2 [47] and *OpenVLA-7B-LIBERO-Long* fine-tuned on LIBERO-Long. During transfer, **no information about victim models** is used, including weights, architecture details beyond public model names, fine-tuning datasets, recipes, or hyperparameters. Specifically, we select *OpenVLA-ofit* [24] and  $\pi$  series [5, 6] models as victim models. Built on OpenVLA, *OpenVLA-ofit* introduces an optimized fine-tuning recipe that notably improves success rates (from 76.5% to 97.1%) and delivers  $\sim 26\times$  throughput. To stress cross-recipe and cross-task vulnerability, we test on four variants fine-tuned on four distinct LIBERO task suites, as well as a multi-suite model trained jointly on all four (*OpenVLA-ofit-w*). The  $\pi$  family differs fundamentally from OpenVLA in backbone choice,

Table 1. Task success rate (%) when transfer from the surrogate OpenVLA-7B to the victim OpenVLA-ofw on LIBERO.

objective	Simulated					Physical				
	spatial	object	goal	long	avg.	spatial	object	goal	long	avg.
Benign	99	99	98	97	98.25	99	99	98	97	98.25
UMA <sub>1</sub>	25	86	40	31	45.50	83	89	76	73	80.25
UMA <sub>1-3</sub>	46	88	38	39	52.75	90	87	83	81	85.25
UADA <sub>1</sub>	35	82	27	21	41.25	71	90	57	74	73.00
UADA <sub>1-3</sub>	37	74	21	33	41.25	65	88	46	61	65.00
TMA <sub>1</sub>	69	89	58	61	69.25	78	92	74	83	81.75
TMA <sub>7</sub>	47	78	47	34	51.50	90	96	89	90	91.25
Our	<b>7</b>	<b>0</b>	<b>10</b>	<b>6</b>	<b>5.75</b>	<b>26</b>	<b>53</b>	<b>54</b>	<b>28</b>	<b>40.25</b>

Table 2. Task success rate (%) when transfer from the surrogate OpenVLA-7B to the victim OpenVLA-of on LIBERO.

objective	Simulated					Physical				
	spatial	object	goal	long	avg.	spatial	object	goal	long	avg.
Benign	98	98	98	94	97.00	98	98	98	94	97.00
UMA <sub>1</sub>	79	95	69	3	61.50	96	90	90	83	89.75
UMA <sub>1-3</sub>	90	93	57	3	60.75	96	81	93	80	87.50
UADA <sub>1</sub>	94	86	61	3	61.00	92	96	79	84	87.75
UADA <sub>1-3</sub>	87	90	64	4	61.25	92	95	61	90	84.50
TMA <sub>1</sub>	99	93	81	25	74.50	98	92	84	86	90.00
TMA <sub>7</sub>	83	93	75	62	78.25	97	88	89	87	90.25
Our	<b>66</b>	<b>43</b>	<b>62</b>	<b>3</b>	<b>43.50</b>	<b>69</b>	<b>74</b>	<b>76</b>	<b>27</b>	<b>61.50</b>

pretraining/fine-tuning data, and training strategy, making transfer substantially harder. We therefore assess black-box transfer on  $\pi_0$  [6], which provide a stringent test of model-agnostic patch behavior across heterogeneous VLA designs.

**Implementation & Evaluation Details.** We evaluate on the LIBERO benchmark [31]. Each suite contains 10 tasks, and each task is attempted in 10 independent trials, yielding 100 rollouts per suite, following [6]. Consistent with [49], patch placement sites are predetermined for each suite to avoid occluding objects in the test scenes. More implementation details can be found in the *Appendix B*. Regarding the evaluation metric, we adopt the concept of Success Rate (SR) introduced in LIBERO [31] across all setting.

## 4.2. Main Results

We first evaluate the white-box performance of our patches, where the victim model is identical to the surrogate. The results in *Appendix C* demonstrate that our method achieves strong white-box attack capability. For the *OpenVLA-7B* [23] to *OpenVLA-ofw* [24] transfer experiment, Tab. 1 shows that our patch objective induces the strongest degradation of task success rates. In the simulated setting, the clean policy succeeds on 98.25% of tasks on average, while our method reduces the success rate to only 5.75%, corresponding to more than a 92% point drop. Existing objectives such as UMA, UADA, and TMA do transfer to

the victim but remain much less destructive: their average success rates stay between 41.25% and 69.25%, and they leave certain categories almost intact, for example, object-centric tasks still above 74% success for UMA and UADA. In contrast, our patch almost completely disables the policy across all four task types. The right block of Tab. 1 reports the attack results under physical setting. A similar trend appears: all baselines still retain high average success (65.00%-91.25%), whereas our method again yields the lowest success rate of 40.25%. This indicates that our patch objective not only transfers more effectively to the simulated environments, but also produces substantially stronger degradation under the physical environment, establishing a consistently harder universal patch baseline across both settings.

Beyond the transfer from *OpenVLA-7B* to *OpenVLA-ofw*, we further evaluate transfer to four different *OpenVLA-of* variants that are separately fine-tuned on different LIBERO task suites, creating a larger distribution and policy gap from the surrogate. Tab. 2 shows that our objective still achieves consistently stronger transfer than all baselines across both simulated and physical setups, highlighting the effectiveness of our design. **Additional transfer results**, including attacks transferred to  $\pi_0$ , are in *Appendix D*, and show that our methods still enhance attacks in the most challenging case of transferring to entirely different

Table 3. Ablation for transfer to openvla-oft under physical setting.

objective	spatial	object	goal	long	avg.
Our	69	74	76	27	61.50
w/o RUPA	70	75	71	33	62.25
w/o PAD	68	67	77	38	62.50
w/o PSM	69	72	81	32	63.50
w/o $\mathcal{J}_{tr}$	90	86	94	73	85.75
w/o $\mathcal{L}_{con}$	93	63	79	48	70.75
w/o $\mathcal{L}_1$	74	74	77	31	64.00

Table 4. Ablation on text-probe phrasing for transfer to openvla-oft in the physical setting.

objective	spatial	object	goal	long	avg.
Our	69	74	76	27	61.50
Action	76	67	94	48	71.25
Direction	72	75	78	75	75.00

VLA.

### 4.3. Ablation Study

**Impact of Each Design.** Tab. 3 further validates the role of each component in our objective. Dropping any single module (RUPA, PAD, or PSM) consistently weakens the attack, reflected by higher average success rates than the full model. The most severe degradation appears in the *w/o*  $\mathcal{J}_{tr}$  variant, where the average success rate jumps to 85.75%, close to the benign and baseline levels. Since  $\mathcal{J}_{tr}$  jointly contains both  $\mathcal{L}_1$  and  $\mathcal{L}_{con}$  (i.e., it removes the entire first-stage feature-space minimization), this indicates that our feature-space  $\ell_1$  and contrastive misalignment objectives, together with the RUPA designs, are essential for strong transfer. Moreover, the impact of  $\mathcal{L}_{con}$  is noticeably larger than that of  $\mathcal{L}_1$ . By Prop. 1 and Cor. 1,  $\mathcal{L}_1$  is a distance-based term that mainly controls the *magnitude* of the surrogate deviation, whereas  $\mathcal{L}_{con}$ , built on cosine similarity, focuses on feature angles and thus shapes the *direction* of the displacement. Consequently, even without  $\mathcal{L}_1$ ,  $\mathcal{L}_{con}$  can still drive patched features away from their clean anchors along transferable directions, so the attack remains relatively strong.

**Impact of Text Probes.** Tab. 4 analyzes how text-probe phrasing influences transfer in the physical setting. We compare our default probes, which jointly encode both action and spatial direction, against two reduced variants: **Action** probes that only include verbs (e.g., “put”, “pick up”, “place”, “turn on”, “push”, “open”, “close”) and **Direction** probes that only contain spatial words (e.g., “left”, “right”, “bottom”, “back”, “middle”, “top”, “front”). Using action-only or direction-only probes markedly weakens the attack: the average success rate increases to 71.25% and 75.00%, respectively, compared to 61.5% with our design. This

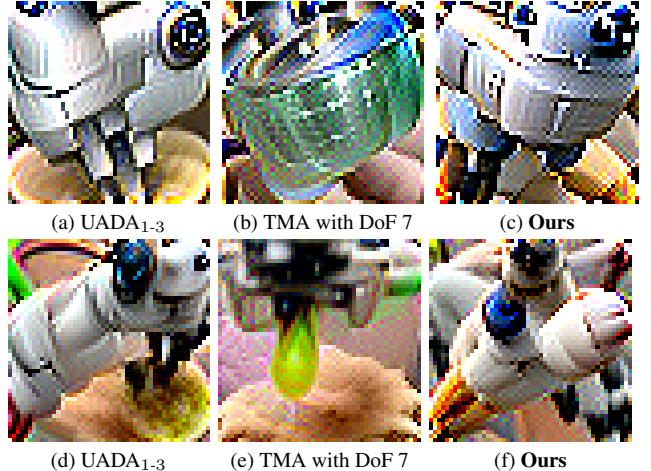


Figure 2. **Patch visualization and comparison.** The first row is trained in a simulated setting, and the second row is trained in a physical setting.

suggests that jointly encoding action and directional cues produces text queries that more closely match the policy’s action-relevant channels, thereby enabling more effective cross-model transfer. Ablation study of more specific parameters can be found in *Appendix E*.

### 4.4. Patch Pattern Analysis

As shown in Fig. 2, we can see that baseline end-to-end methods [49] produce scene-tied patterns: UADA yields textures that closely resemble the robot gripper in both simulation and physical settings (Fig. 2a and 2d), while TMA generates more abstract yet surrogate-specific shapes (Fig. 2b and 2e). These behaviors indicate overfitting to object/embodyment cues, which hampers cross-model and cross-setting transfer. In contrast, our universal transferable patch (Fig. 2c and 2f) is learned in feature space to perturb higher-level, model-agnostic representations shared across VLAs. By jointly optimizing feature-space, attention, and semantic objectives, our patch combines the strengths of prior designs, avoids object mimicry, and yields a universal patch that reliably transfers across tasks, embodiments, and environments, resulting in stronger black-box transfer.

## 5. Conclusion

In this paper, we present the first study of universal, transferable patch attacks on VLA-driven robots and introduce UPA-RFAS, a unified framework that couples an  $\ell_1$  feature deviation with repulsive contrastive alignment to steer perturbations into model-agnostic, high-transfer directions. UPA-RFAS integrates a robustness-augmented patch optimization and two VLA-specific losses, Patch Attention Dominance and Patch Semantic Misalignment, which achieve strong black-box transfer across models, tasks, and



sim-to-real settings, revealing a practical patch-based threat and a solid baseline for future defenses.

## References

- [1] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016. 3
- [2] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pages 274–283. PMLR, 2018. 3
- [3] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024. 2
- [4] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025. 2
- [5] Kevin Black et al.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization. 2024. 1, 6
- [6] Kevin Black et al.  $\pi_0$ : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 1, 2, 6, 7
- [7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 1, 2
- [8] Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017. 1, 2
- [9] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017. 2
- [10] Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: elastic-net attacks to deep neural networks via adversarial examples. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 1
- [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmlR, 2020. 1, 4
- [12] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *CVPR*, 2018. 2
- [13] Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *CVPR*, 2019. 2
- [14] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Wei, Yongdae Bo, Amir Rahmati, Dawn Song, Patrick Traynor, Atul Prakash, and Tadayoshi Kohno. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [15] Senyu Fei, Siyin Wang, Junhao Shi, Zihao Dai, Jikun Cai, Pengfang Qian, Li Ji, Xinzhe He, Shiduo Zhang, Zhaoye Fei, et al. Libero-plus: In-depth robustness analysis of vision-language-action models. *arXiv preprint arXiv:2510.13626*, 2025. 1
- [16] Aditya Ganeshan, Vivek BS, and R Venkatesh Babu. Fda: Feature disruptive attack. In *ICCV*, 2019. 2
- [17] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015. 2
- [18] Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-Jen Wang, Yucheng Hu, and Jianyu Chen. Improving vision-language-action model with online reinforcement learning. *arXiv preprint arXiv:2501.16664*, 2025. 2
- [19] Hao Huang, Ziyang Chen, Huanran Chen, Yongtao Wang, and Kevin Zhang. T-sea: Transfer-based self-ensemble attack on object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20514–20523, 2023. 2
- [20] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023. 3
- [21] Xiaojun Jia, Sensen Gao, Simeng Qin, Tianyu Pang, Chao Du, Yihao Huang, Xinfeng Li, Yiming Li, Bo Li, and Yang Liu. Adversarial attacks against closed-source mllms via feature optimal alignment. *arXiv preprint arXiv:2505.21494*, 2025. 1
- [22] Haydn T Jones, Jacob M Springer, Garrett T Kenyon, and Juston S Moore. If you’ve trained one you’ve trained them all: inter-architecture similarity increases with robustness. In *Uncertainty in Artificial Intelligence*, pages 928–937. PMLR, 2022. 4
- [23] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 1, 2, 3, 6, 7
- [24] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025. 1, 6, 7
- [25] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pages 3519–3529. PMIR, 2019. 3
- [26] Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *Artificial intelligence safety and security*, pages 99–112. Chapman and Hall/CRC, 2018. 2
- [27] Qixiu Li, Yaobo Liang, Zeyu Wang, Lin Luo, Xi Chen, Mozheng Liao, Fangyun Wei, Yu Deng, Sicheng Xu, Yizhong Zhang, et al. Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation. *arXiv preprint arXiv:2411.19650*, 2024. 2

- [28] Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18061–18070, 2024. 2
- [29] Xiao Li, Yiming Zhu, Yifan Huang, Wei Zhang, Yingzhe He, Jie Shi, and Xiaolin Hu. Pbcats: Patch-based composite adversarial training against physically realizable attacks on object detection. *arXiv preprint arXiv:2506.23581*, 2025. 1
- [30] Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E Hopcroft. Nesterov accelerated gradient and scale invariance for adversarial attacks. *arXiv preprint arXiv:1908.06281*, 2019. 2
- [31] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36:44776–44791, 2023. 1, 6, 7
- [32] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *ICLR*, 2017. 2
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [34] Guanxing Lu, Wenkai Guo, Chubin Zhang, Yuheng Zhou, Haonan Jiang, Zifeng Gao, Yansong Tang, and Ziwei Wang. Vla-rl: Towards masterful and general robotic manipulation with scalable reinforcement learning. *arXiv preprint arXiv:2505.18719*, 2025. 2
- [35] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017. 6
- [36] Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. *Advances in neural information processing systems*, 31, 2018. 3
- [37] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 3
- [38] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025. 2
- [39] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. *Advances in neural information processing systems*, 30, 2017. 1, 3
- [40] Alexander Robey, Zachary Ravichandran, Vijay Kumar, Hamed Hassani, and George J Pappas. Jailbreaking llm-controlled robots. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11948–11956. IEEE, 2025. 1
- [41] Prashant Shekhar, Bidur Devkota, Dumindu Samaraweera, Laxima Niure Kandel, and Manoj Babu. Do adversarial patches generalize? attack transferability study across real-time segmentation models in autonomous vehicles. In *2025 IEEE Security and Privacy Workshops (SPW)*, pages 322–328. IEEE, 2025. 2
- [42] Jacob Springer, Melanie Mitchell, and Garrett Kenyon. A little robustness goes a long way: Leveraging robust features for targeted transfer attacks. *Advances in Neural Information Processing Systems*, 34:9759–9773, 2021. 4
- [43] Andreas Steiner, André Susano Pinto, Michael Tschannen, Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey Gritsenko, Matthias Minderer, Anthony Sherbondy, Shangbang Long, et al. Paligemma 2: A family of versatile vlms for transfer. *arXiv preprint arXiv:2412.03555*, 2024. 2
- [44] Shuhan Tan, Kairan Dou, Yue Zhao, and Philipp Krähenbühl. Interactive post-training for vision-language-action models. *arXiv preprint arXiv:2505.17016*, 2025. 2
- [45] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 3
- [46] Florian Tramer. Detecting adversarial examples is (nearly) as hard as classifying them. In *International Conference on Machine Learning*, pages 21692–21702. PMLR, 2022. 3
- [47] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, 2023. 1, 6
- [48] Kunyu Wang, Xuanran He, Wenxuan Wang, and Xiaosen Wang. Boosting Adversarial Transferability by Block Shuffle and Rotation. In *CVPR*, 2024. 2
- [49] Taowen Wang, Cheng Han, James Liang, Wenhao Yang, Dongfang Liu, Luna Xinyu Zhang, Qifan Wang, Jiebo Luo, and Ruixiang Tang. Exploring the adversarial vulnerabilities of vision-language-action models in robotics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6948–6958, 2025. 1, 3, 6, 7, 8
- [50] Xiaosen Wang and Kun He. Enhancing the transferability of adversarial attacks through variance tuning. In *CVPR*, 2021. 2
- [51] Xiaosen Wang, Xuanran He, Jingdong Wang, and Kun He. Admix: Enhancing the transferability of adversarial attacks. In *ICCV*, 2021. 2
- [52] Zhibo Wang, Hengchang Guo, Zhifei Zhang, Wenxin Liu, Zhan Qin, and Kui Ren. Feature importance-aware transferable adversarial attacks. In *ICCV*, 2021. 2
- [53] Junjie Wen, Yichen Zhu, Jinming Li, Minjie Zhu, Zhibin Tang, Kun Wu, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *IEEE Robotics and Automation Letters*, 2025. 2
- [54] Junjie Wen, Yichen Zhu, Minjie Zhu, Zhibin Tang, Jinming Li, Zhongyi Zhou, Xiaoyu Liu, Chaomin Shen, Yaxin

- Peng, and Feifei Feng. Diffusionvla: Scaling robot foundation models via unified diffusion and autoregression. In *Forty-second International Conference on Machine Learning*, 2025. [2](#)
- [55] Zihao Xiao, Xianfeng Gao, Chilin Fu, Yinpeng Dong, Wei Gao, Xiaolu Zhang, Jun Zhou, and Jun Zhu. Improving transferability of adversarial patches on face recognition with generative models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11845–11854, 2021. [2](#)
- [56] Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L Yuille. Improving transferability of adversarial examples with input diversity. In *CVPR*, 2019. [2](#)
- [57] Haochuan Xu, Yun Sing Koh, Shuhuai Huang, Zirun Zhou, Di Wang, Jun Sakuma, and Jingfeng Zhang. Model-agnostic adversarial attack and defense for vision-language-action models. *arXiv preprint arXiv:2510.13237*, 2025. [1](#)
- [58] Kaidi Xu, Gaoyuan Zhang, Sijia Liu, Quanfu Fan, Mengshu Sun, Hongge Chen, Pin-Yu Chen, Yanzhi Wang, and Xue Lin. Adversarial t-shirt! evading person detectors in a physical world. In *European conference on computer vision*, pages 665–681. Springer, 2020. [2](#)
- [59] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023. [2](#), [3](#)
- [60] Hangtao Zhang, Chenyu Zhu, Xianlong Wang, Ziqi Zhou, Changgan Yin, Minghui Li, Lulu Xue, Yichen Wang, Shengshan Hu, Aishan Liu, et al. Badrobot: Jailbreaking embodied llms in the physical world. *arXiv preprint arXiv:2407.20242*, 2024. [1](#)
- [61] Jianping Zhang, Weibin Wu, Jen-tse Huang, Yizhan Huang, Wenxuan Wang, Yuxin Su, and Michael R Lyu. Improving adversarial transferability via neuron attribution-based attacks. In *CVPR*, 2022. [2](#)
- [62] Yaoyuan Zhang, Yu-an Tan, Tian Chen, Xinrui Liu, Quanxin Zhang, and Yuanzhang Li. Enhancing the transferability of adversarial examples with random patch. In *IJCAI*, 2022. [2](#)
- [63] Xueyang Zhou, Guiyao Tie, Guowen Zhang, Hechang Wang, Pan Zhou, and Lichao Sun. Badvla: Towards backdoor attacks on vision-language-action models via objective-decoupled optimization. *arXiv preprint arXiv:2505.16640*, 2025. [1](#)
- [64] Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiayi Cui, Hongfa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, et al. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. *arXiv preprint arXiv:2310.01852*, 2023. [2](#)
- [65] Hegui Zhu, Yuchen Ren, Xiaoyan Sui, Lianping Yang, and Wuming Jiang. Boosting adversarial transferability via gradient relevance attack. In *ICCV*, 2023. [2](#)
- [66] Rongyi Zhu, Zeliang Zhang, Susan Liang, Zhuo Liu, and Chenliang Xu. Learning to transform dynamically for better adversarial transferability. In *CVPR*, 2024. [2](#)
- [67] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023. [1](#), [2](#)