# Eva-VLA: Evaluating Vision-Language-Action Models' Robustness Under Real-World Physical Variations

Hanqing Liu[1,2,3], Jiahuan Long[1,2,3], Junqi Wu[1,2,3], Jiacheng Hou[2,3], Huili Tang[1,2,3],
Tingsong Jiang[2,3*], Weien Zhou[2,3], Wen Yao[2,3*]

[1]MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
[2]Defense Innovation Institute, Chinese Academy of Military Science
[3]Intelligent Game and Decision Laboratory

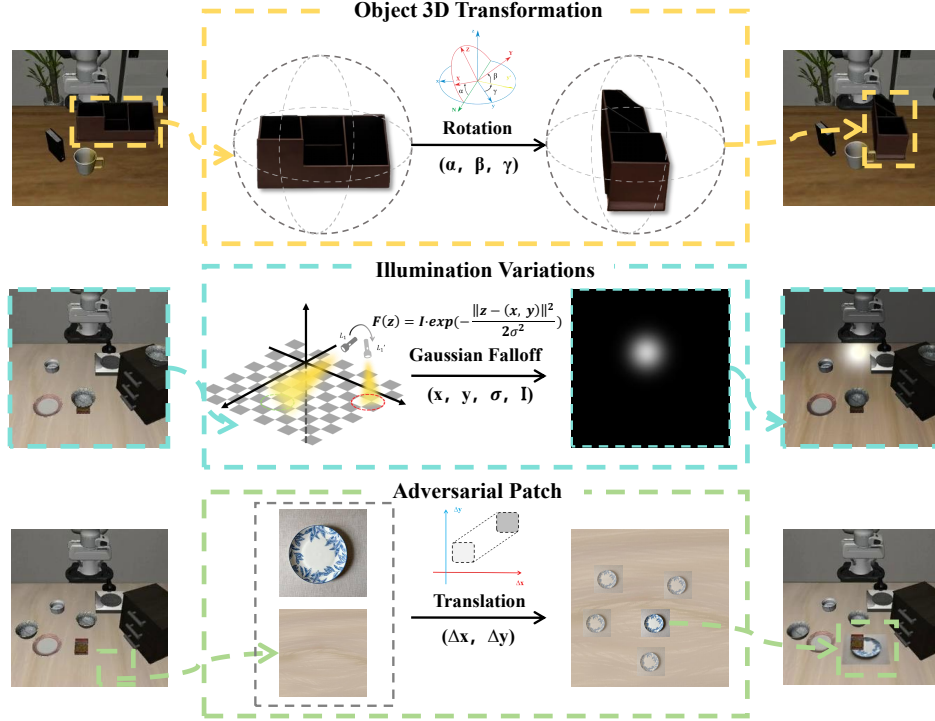hanqingliu@sjtu.edu.cn,tingsong@pku.edu.cn,wendy0782@126.com

Fig. 1: **Visualization of three categories of physical variations.** Object 3D transformations through rotation parameters $(\alpha, \beta, \gamma)$ that alter object 3D poses in the scene **(top)**. Illumination variations modeled as Gaussian falloff functions with parameters $(x, y, \sigma, I)$ controlling illumination position, radius, and intensity **(middle)**. Adversarial patches with translation parameters $(\Delta x, \Delta y)$ that introduce visual disruptions at critical locations in the scene **(bottom)**.

*Abstract*— Vision-Language-Action (VLA) models have emerged as promising solutions for robotic manipulation, yet their robustness to real-world physical variations remains critically underexplored. To bridge this gap, we propose Eva-VLA, the first unified framework that systematically evaluates the robustness of VLA models by transforming discrete physical variations into continuous optimization problems. However, comprehensively assessing VLA robustness presents two key challenges: (1) how to systematically characterize diverse physical variations encountered in real-world deployments while maintaining evaluation reproducibility, and (2) how to discover worst-case scenarios without prohibitive real-world data collection costs efficiently. To address the first challenge, we decompose real-world variations into three critical domains: object 3D transformations that affect spatial reasoning, illumination variations that challenge visual perception, and adversarial patches that disrupt scene understanding. For the second challenge, we introduce a continuous black-box optimization framework that transforms discrete physical variations into parameter optimization, enabling systematic exploration of worst-case scenarios. Extensive experiments on state-of-the-art OpenVLA models across multiple benchmarks reveal alarming vulnerabilities: all variation types trigger failure rates exceeding **60%**, with object transformations causing up to **97.8%** failure in long-horizon tasks. Our findings expose critical gaps between controlled laboratory success and unpredictable deployment readiness, while the Eva-VLA framework provides a practical pathway for hardening VLA-based robotic manipulation models against real-world deployment challenges.

* Corresponding author.

## I. INTRODUCTION

Vision-Language-Action (VLA) models represent a paradigm shift in robotic manipulation, integrating visual perception, language understanding, and action generation into unified end-to-end systems [1]. Recent deployments across manufacturing [2], healthcare [3], and service robotics [4], [5] demonstrate their transformative potential. However, in real-world deployments, VLA models inevitably face challenging physical variations, such as spatial transformations, illumination variations, and visual disruptions, which can dramatically alter robot behavior without being immediately detectable, posing significant safety risks. Therefore, it is crucial to investigate VLA robustness across various physical conditions systematically.

Existing research has explored the robustness of VLA-based robotic systems through approaches like adversarial patches [6], which generate localized perturbations via gradient-based white-box attacks to achieve visual interference. However, these methods suffer from critical limitations: they violate physical plausibility constraints and fail to capture the rich spectrum of real-world physical variations. Moreover, their reliance on gradient access restricts applicability to black-box deployment scenarios. Based on these limitations, we aim to generate more diverse and realistic physical variations for comprehensively evaluating VLA robustness, while two key challenges must be addressed: **(1)** *How to systematically characterize diverse physical variations encountered in real-world deployments while maintaining evaluation reproducibility?* **(2)** *How to discover worst-case scenarios without prohibitive real-world data collection costs efficiently?*

To address these challenges, we propose *Eva-VLA*, a unified framework for evaluating vision-language-action models' robustness. Our key innovation lies in transforming discrete physical variations into continuous optimization problems. First, as shown in Fig. 1, we decompose real-world variations into three distinct domains: object 3D transformations parameterized with rotation angles($\alpha$, $\beta$, $\gamma$), illumination variations defined by point light parameters including position($x$, $y$), radius($\sigma$), intensity($I$), and adversarial patch placement specified by ($\Delta x$, $\Delta y$). This parameterization enables systematic exploration of the variation space while maintaining physical plausibility through explicit constraints. Second, to overcome the black-box nature of VLA models and non-differentiable simulation environments, we employ Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [7], a gradient-free optimization algorithm, to efficiently discover worst-case scenarios by iteratively optimizing physical variations parameters. This approach enables comprehensive vulnerability assessment without requiring model gradients or expensive real-world data collection.

Our main contributions are as follows: ❶ To the best of our knowledge, we are the first to decompose real-world physical variations into three key domains—object 3D transformation, illumination variations, and adversarial patches—enabling a comprehensive evaluation of VLA robustness under these physical variations. ❷ We propose *Eva-VLA*, a novel framework that transforms discrete physical variations into continuous parameter optimization. By leveraging a simulator environment that allows us to reset to the same conditions, we ensure the repeatability and reliability of the evaluation process, which enables efficient exploration of worst-case scenarios without the need for expensive real-world data collection. ❸ Through extensive experiments on state-of-the-art model OpenVLA across multiple benchmarks, we expose significant fragility in current VLA systems, with failure rates exceeding 60% across all variation categories, with object transformations causing up to 97.8% failure in long-horizon tasks. These findings provide crucial insights for developing more robust VLA architectures and underscore the urgent need for improved robustness training methodologies.

## II. RELATED WORK

### A. Vision-Language-Action Models

Recent studies [8], [9], [4], [2], [10], [11], [12], [13] have explored extending pre-trained Vision-Language Models (VLMs) to generate robot actions, thereby contributing to the development of generalist robotic policies. RT-2 [9] pioneered this approach by fine-tuning PaLI-X [14] with robot actions discretized into 256 bins as tokens. Building on this foundation, OpenVLA [4] applies a similar discretization but fine-tunes Prismatic VLM [15] specifically on the OXE dataset [16], which contains data from 22 robot embodiments. The OpenVLA architecture has spawned several improvements. For example, CogACT [2] incorporates a diffusion-based action module for better continuous action modeling. TraceVLA [10] introduces visual trace prompting for enhanced spatial-temporal awareness, achieving 3.5x performance gains on real robot tasks. SpatialVLA [17] adds 3D position encoding and adaptive action grids for improved spatial understanding. Notably, OpenVLA-OFT [18] optimizes fine-tuning through parallel decoding and action chunking, achieving 26x speedup. We select OpenVLA series as our target models due to its comprehensive, open-source framework and widespread community adoption, which ensures reproducible comparisons.

### B. Adversarial Attacks in Robotic Systems

Adversarial attacks play a crucial role in assessing the vulnerabilities of machine learning models, particularly in robotics, where models operate in dynamic, real-world environments. Traditional gradient-based, pixel-level attacks [19], [20], [21], [22], [23] exploit model gradients to calculate malicious perturbations, achieving high success rates in controlled digital settings. However, when applied to real-world scenarios, these attacks often face significant challenges due to the complexity and variability of physical environments. For physical-world attacks, patch-based methods [24], [25] have emerged as a practical solution, with recent works examining illumination variations [26] and viewpoint variations [27] impact model robustness. Since VLA models directly output control signals for robotic actions through end-to-end learning, attacks targeting the visual input can
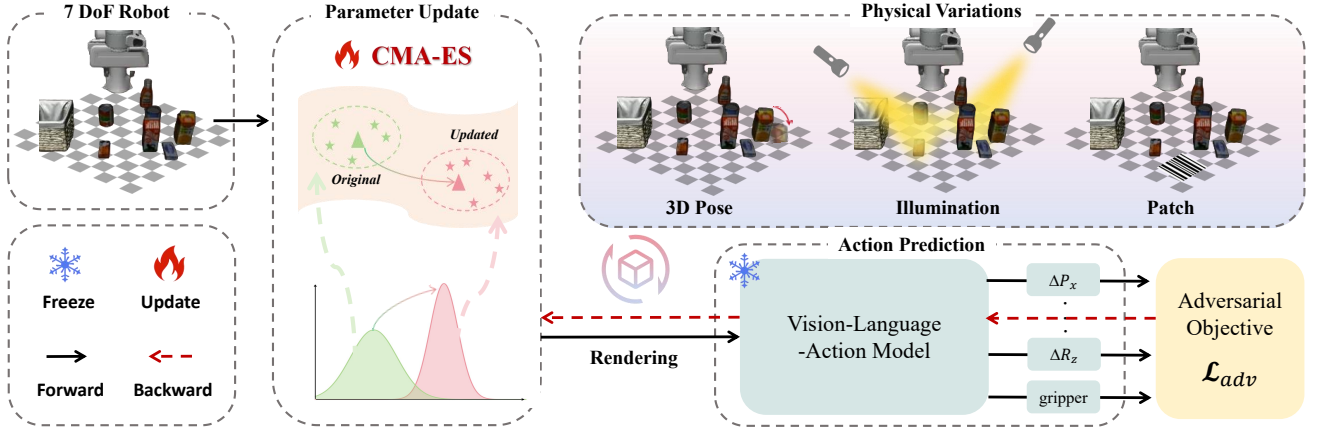
Fig. 2: **Overview of the Proposed Eva-VLA Framework.** To capture worst-case physical variations, discrete transformations in three critical domains are parameterized and their distributions optimized through a query-based method, maximizing prediction errors of vision-language-action models to reduce task success rates.

significantly affect the final action outputs. In this context, Wang et al. [6] explored the vulnerabilities of VLA systems by placing a colored patch within the camera's field of view. While this approach led to a decrease in task success rates, it relied on white-box optimization and produced patches that were easily detectable by the human eye. We present the first comprehensive study of VLA robustness under physical transformations—3D poses, illumination, and adversarial patches—addressing real-world deployment challenges.

## III. METHOD

In this section, we first provide an overview of VLA models and their underlying principles. Then, we present our Eva-VLA framework (illustrated in Fig. 2), detailing the parameterization of three categories of physical variations and the formulation of our adversarial objectives. Finally, we describe the optimization algorithm that efficiently discovers adversarial distributions over these physical variations.

### A. Preliminaries

Vision-Language-Action (VLA) models are built on large language models integrated with visual encoders, enabling end-to-end task execution by jointly processing visual perception, linguistic instructions, and action generation. Formally, a VLA model can be defined as: $F : \mathcal{V} \times \mathcal{L} \to \mathcal{A}$ where $\mathcal{V}$ represents the visual observation space, $\mathcal{L}$ denotes the language instruction space, and $\mathcal{A}$ is the action output space. In this work, we focus on a 7 degrees of freedom (DoF) robotic manipulator [28], where the action space encodes both end-effector motion and gripper control. The output action vector can be specified as:

$$A = [\Delta P_x, \Delta P_y, \Delta P_z, \Delta R_x, \Delta R_y, \Delta R_z, gripper], \quad (1)$$

where $\Delta P_{x,y,z} \in \mathbb{R}^3$ and $\Delta R_{x,y,z} \in \mathbb{R}^3$ represent the positional and rotational changes along the Cartesian axes, and $gripper \in \mathbb{R}$ controls the gripper states [4].

### B. Eva-VLA Framework

To comprehensively evaluate the VLA model's robustness, we introduce the first unified framework Eva-VLA in Fig. 2,

which parameterizes three categories of physical variations that span different aspects of visual variations encountered in real-world scenarios. ❶ **Parametrization of 3D Transformation.** In this section, we primarily focus on rigid 3D transformations of important objects in the scene, such as the rotation, since it reflects the most typical 3D changes observed in real-world environments. Formally, we define a 3-dimensional vector $\boldsymbol{\Theta} = \{\alpha, \beta, \gamma\}$ to uniquely parameterize any arbitrary transformation, which denotes the Tait-Bryan angles (yaw, pitch, roll) in the $Z$-$Y$-$X$ sequence. In order to comply with the physical laws of the space in which the target is placed, we need to restrict its transformation range to a bounded range of $[\boldsymbol{\Theta}_{\min}, \boldsymbol{\Theta}_{\max}]$. ❷ **Parametrization of Illumination Variations.** In this section, we primarily focus on the illumination variations in real-world environments. Formally, we define a 4-dimensional vector $\boldsymbol{\Lambda} = \{x, y, \sigma, I\}$ to parameterize the illumination environment. Then, to simulate realistic light behavior, we apply a Gaussian falloff function to model the spatial distribution of the light source, which can be represented as:

$$L(z) = I \cdot \exp\left(-\frac{\|z - (x,y)\|^2}{2\sigma^2}\right), \quad (2)$$

where $z$ represents any point in the scene. The parameters $(x, y)$, $\sigma$, and $I$ from $\boldsymbol{\Lambda}$ directly control the center position, brightness, and spread of a point light source, respectively. ❸ **Parametrization of Adversarial Patch.** Rather than optimizing patch textures [6], we employ natural images (e.g., barcodes, QR codes, everyday images) and optimize their placement on the tabletop surface. This approach ensures physical realizability through standard UV mapping while maximizing perceptual impact, as the table occupies a significant portion of the robot's visual field during manipulation. Formally, we define a 2-dimensional vector: $\boldsymbol{\phi} = \{x, y\} \in [\boldsymbol{\phi}_{min}, \boldsymbol{\phi}_{max}]$, where $(x, y)$ denotes the position of the adversarial patch in the table texture. This constraint ensures patches remain within the robot's primary workspace, reducing optimization queries while maintaining adversarial effectiveness.

Then, we formulate the adversarial objective to systematically evaluate how these variations impact VLA model behavior. Specifically, our adversarial objective $\mathcal{L}_{adv}$ is to cause the 7-Dof robot arm to output incorrect action prediction vectors, which can be expressed as:

$$\mathcal{L}_{adv} = -\sum_{i=1}^{N} cos(A_{clean}^i, A_{adv}^i), \quad (3)$$

where $N$ denotes the total number of action vector, $A^i$ denotes the i-th action vector, which can be represented as:

$$A_{adv}^i = F(T(X^i, \boldsymbol{\mathcal{T}}), Y), \quad (4)$$

where $F$ represents the specified VLA models, $T(\cdot)$ is a transformation function, $X$ denotes the input image, $Y$ denotes the input instruction. Thus, the objective of Eva-VLA is to discover an optimal distribution $p^*(\boldsymbol{\mathcal{T}})$ over transformation parameters $\boldsymbol{\mathcal{T}} \in \{\boldsymbol{\Theta}, \boldsymbol{\Lambda}, \boldsymbol{\phi}\}$. The optimization seeks to ensure that parameters sampled from this distribution maximize adversarial impact on VLA models. This can be formulated as:

$$p^*(\boldsymbol{\mathcal{T}}) = \arg\max_{p(\boldsymbol{\mathcal{T}})} \mathbb{E}_{\boldsymbol{\mathcal{T}} \sim p(\boldsymbol{\mathcal{T}})}[\mathcal{L}_{adv}(A_{clean}, A_{adv})]. \quad (5)$$

### C. Query-Based Optimizing Algorithm

In this section, we present our physical variation optimization algorithm, which aims to identify optimal transformation configurations $\boldsymbol{\mathcal{T}}^*$. Rather than optimizing for a single configuration, we explore a distribution over configurations, offering several key advantages: **1)** It enables diverse exploration of the configuration space, helping identify regions where optimal adversarial conditions exist; **2)** It facilitates faster convergence to regions containing effective configurations while reducing the risk of local optima; **3)** It enhances robustness by avoiding overfitting to specific solutions, ensuring better generalization across different VLA models.

To unify the optimization of different physical variations—3D transformations, illumination variations, and adversarial patch—we parameterize each discrete variation type as a continuous distribution. This parameterization transforms the discrete search problem into a continuous optimization problem, enabling efficient exploration of the adversarial space. Specifically, each variation type $\boldsymbol{\mathcal{T}} \in \{\boldsymbol{\Theta}, \boldsymbol{\Lambda}, \boldsymbol{\phi}\}$ is modeled as a multivariate Gaussian distribution: $\boldsymbol{\mathcal{T}} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathcal{T}}, \boldsymbol{\Sigma}_{\mathcal{T}}^2 C_{\mathcal{T}})$, where $\boldsymbol{\mu}_{\mathcal{T}}$ represents the mean configuration, $\boldsymbol{\Sigma}_{\mathcal{T}}$ controls the exploration scale, and $C_{\mathcal{T}}$ captures parameter correlations. While all three variation types share the same initial distribution structure, they evolve independently during optimization to discover variation-specific adversarial patterns.

We employ Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [7], [29], a gradient-free optimization algorithm particularly well-suited for black-box adversarial optimization. CMA-ES does not require gradient information, making it ideal for querying VLA models that may not provide gradient access. The algorithm adaptively learns the distribution parameters by iteratively sampling, evaluating,

---

**Algorithm 1:** Optimization Algorithm

**Data:** Task dataset $\mathcal{D} = \{(X_j, Y_j, A_j)\}_{j=1}^N$, VLA model $F$, Attack objective $\mathcal{L}_{adv}$, Transformation function $T$, Initial distribution parameters $\boldsymbol{\mu}_0, \boldsymbol{C}_0, \boldsymbol{\Sigma}_0$.

**Result:** Adversarial transformation configuration $\boldsymbol{\mathcal{T}}^*$.

/* Initialization of distribution parameters */
1   $\boldsymbol{\mu}_0 \leftarrow \mathbf{0}, \boldsymbol{C}_0(\Sigma_0)^2 \leftarrow \mathbf{I}$;
2   **while** $t < t_{\max}$ **do**

    /* Step 1: Sample and evaluate */
3     **for** $i = 1 \rightarrow K$ **do**
4        $\boldsymbol{\mathcal{T}}_i^{t+1} \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{C_t}(\Sigma_t)^2)$;
5        $(A_{adv})_i^{t+1} \leftarrow F(T(\boldsymbol{\mathcal{T}}_i^{t+1}, X), Y)$;

      /* Calculate loss value */
6        $\mathcal{L}_i^{t+1} = \mathcal{L}_{adv}((A_{adv})_i^{t+1}, A_{clean})$;

      /* Early Stopping */
7        **if** *should_stop()* **then**
8          break

    /* Step 2: Selection */
9     Sorting $\mathcal{L}_i^{t+1}$ in ascending order;

    /* Step 3: Update */
10    $\boldsymbol{\mu}_{t+1} \leftarrow \sum_{i=1}^K w_i \boldsymbol{\mathcal{T}}_i^{t+1}$;
11    $\boldsymbol{C}_{t+1} \leftarrow$
     $(1 - c_c)\boldsymbol{C}_t + c_c \sum_{i=1}^K w_i (\boldsymbol{\mathcal{T}}_i^{t+1} - \boldsymbol{\mu}_t)(\boldsymbol{\mathcal{T}}_i^{t+1} - \boldsymbol{\mu}_t)^T$;
12    $\boldsymbol{\Sigma}_{t+1} \leftarrow \boldsymbol{\Sigma}_t \cdot \exp\left(c_\sigma\left(\frac{\|\mathbf{p}\|}{\mathbb{E}[\|\mathcal{N}(0,I)\|]} - 1\right)\right)$;

13   $\boldsymbol{\mu}^* \leftarrow \boldsymbol{\mu}_{t_{\max}}, \boldsymbol{C}^* \leftarrow \boldsymbol{C}_{t_{\max}}, \boldsymbol{\Sigma}^* \leftarrow \boldsymbol{\Sigma}_{t_{\max}}$;
14   $\boldsymbol{\mathcal{T}}^* \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{C}^*(\Sigma^*)^2)$;

---

and updating based on the adversarial loss.

To enhance optimization efficiency, we incorporate two key techniques: ❶ **Learning Rate Adaptation (LRA):** This technique dynamically adjusts the step size $\boldsymbol{\Sigma}$ during optimization to balance exploration and exploitation, thereby improving convergence speed. ❷ **Early Stopping Policy:** Monitors the convergence criterion and terminates optimization when improvements plateau, preventing computational waste on marginal gains. The optimization process iterates through three main steps: sampling candidate configurations from the current distribution, evaluating adversarial effectiveness through VLA model queries, and updating distribution parameters based on the most successful candidates. The complete procedure is detailed in Algorithm 1.

## IV. EXPERIMENTS

### A. Implementation Details

We carefully design the boundaries of transformations parameters to ensure both adversarial effectiveness and physical plausibility in real-world deployment.

**Object 3D Transformations.** Due to rigid body physics and gravitational constraints, objects in tabletop manipulation must maintain stable contact with the surface. Under reasonable physical assumptions, we parameterize rotations exclusively around the vertical axis (z-axis) perpendicular to the table plane, with $\gamma \in [-90, 90]$. This design choice ensures that objects remain grounded on the table surface without floating or violating physical laws. Rotations around other axes are constrained to zero, preventing physically
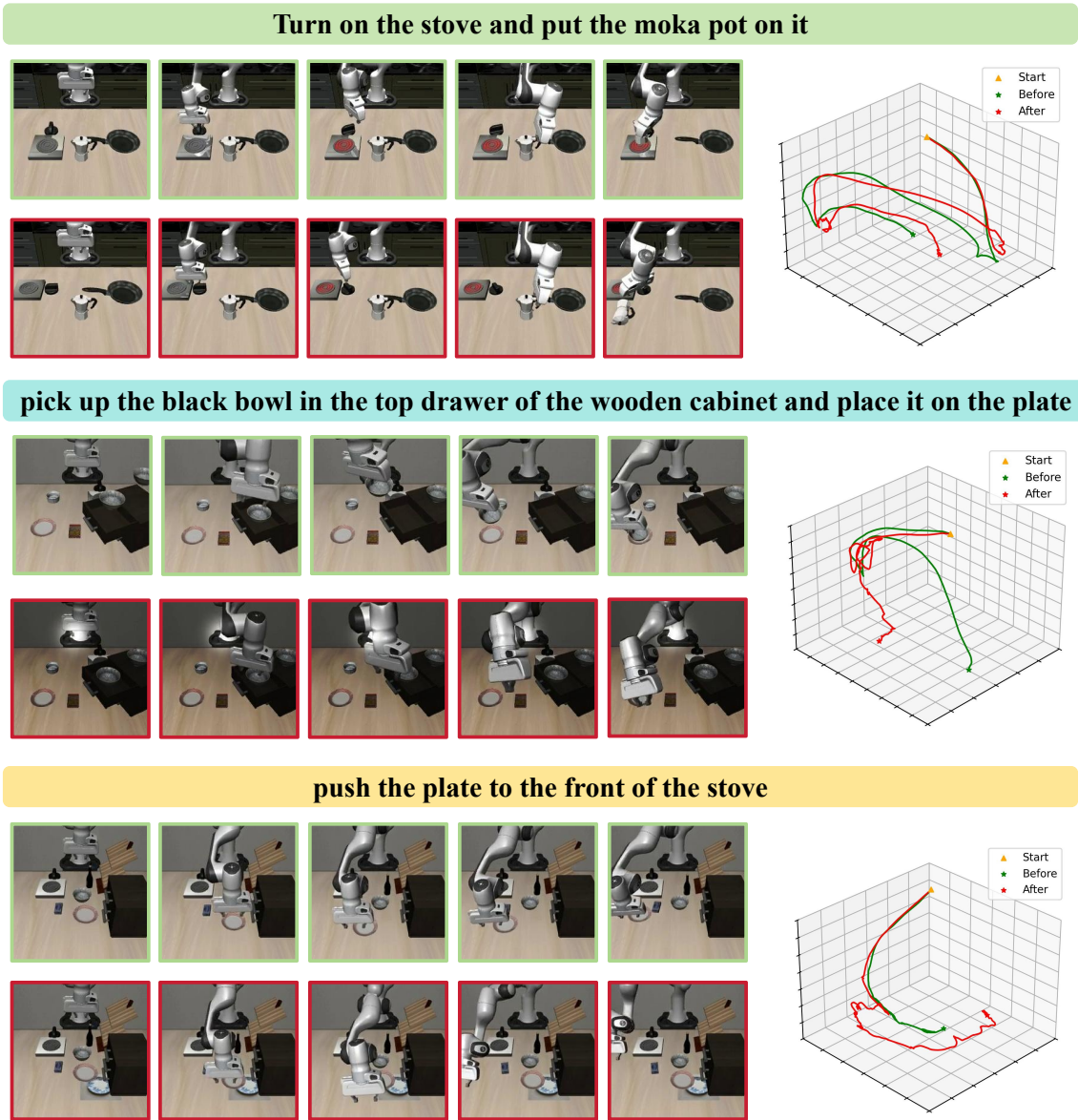
Fig. 3: **Qualitative results under three physical variations on OpenVLA-7B [4] fine-tuned on LIBERO [30].** Three manipulation tasks are shown with original executions (top row) and adversarially perturbed executions (bottom row, highlighted in red) for each task. The 3D trajectory visualizations on the right demonstrate the end-effector paths before (**green**) and after (**red**) applying physical variations, illustrating how object 3D transformations (**top**), illumination variations (**middle**), and adversarial patches (**bottom**) respectively disrupt the robot's motion patterns and lead to task failures.

implausible configurations such as objects suspended in mid-air or penetrating the table surface. This parameterization effectively captures natural pose variations while maintaining strict physical realism.

**Illumination Variations.** For illumination variations, we model point light sources that can be positioned at arbitrary locations within the scene. The light position is parameterized to move freely within the workspace, as even boundary-positioned lights can significantly alter the global illumination distribution. To ensure a fair comparison across experiments, we standardize the light source configuration to a fixed radius of 50 pixels and an intensity of 0.8. These values are chosen to create noticeable but realistic lighting

variations that mirror common indoor lighting conditions encountered in robotic deployment environments. The point light model with Gaussian falloff ensures smooth, physically plausible illumination changes across the scene.

**Adversarial Patch Placement.** To maximize the impact on visual perception, adversarial patches must cover substantial areas within the robot's observation space. We constrain patch placement to the central region of the table texture, specifically:

$$\phi = \{x, y\} \in \left[\frac{W}{3}, \frac{2W}{3}\right] \times \left[\frac{H}{3}, \frac{2H}{3}\right], \quad (6)$$

where $W$ and $H$ denote the table texture dimensions. This

TABLE I: **Main Results.** Failure Rate (FR↑) of OpenVLA and OpenVLA-OFT (fine-tuned) across different tasks within LIBERO suite are reported. * represents the evaluation results of the corresponding model variant on the dataset, $\Theta$ denotes the object 3d transformation, $\Lambda$ denotes the illumination variations, $\phi$ denotes the adversarial patch placements. Number in **bold** indicates the best performance.

| Model | Type | Method | Libero [30] | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | | Spatial* | Object* | Goal* | Long* | Avg |
| OpenVLA (fine-tuned) [4] | − | Clean | $15.3 \pm 0.9\%$ | $11.6 \pm 0.8\%$ | $20.8 \pm 1.0\%$ | $46.3 \pm 1.3\%$ | $23.5 \pm 0.9\%$ |
| | $\Theta$ | Random | $48.2 \pm 3.2\%$ | $52.6 \pm 5.2\%$ | $45.3 \pm 6.8\%$ | $76.8 \pm 11.2\%$ | $55.7 \pm 6.6\%$ |
| | | Best | $\mathbf{71.2 \pm 2.4\%}$ | $\mathbf{78.8 \pm 4.8\%}$ | $\mathbf{82.6 \pm 3.2\%}$ | $\mathbf{97.8 \pm 4.8\%}$ | $\mathbf{82.6 \pm 3.8\%}$ |
| | $\Lambda$ | Random | $28.2 \pm 5.0\%$ | $34.4 \pm 6.4\%$ | $29.6 \pm 4.4\%$ | $60.2 \pm 12.4\%$ | $38.1 \pm 7.1\%$ |
| | | Best | $\mathbf{54.6 \pm 3.6\%}$ | $\mathbf{48.2 \pm 2.8\%}$ | $\mathbf{58.6 \pm 4.2\%}$ | $\mathbf{88.8 \pm 3.4\%}$ | $\mathbf{62.6 \pm 3.5\%}$ |
| | $\phi$ | Random | $22.8 \pm 5.8\%$ | $23.6 \pm 4.6\%$ | $34.2 \pm 7.2\%$ | $52.2 \pm 7.8\%$ | $33.2 \pm 6.2\%$ |
| | | Best | $\mathbf{59.6 \pm 4.4\%}$ | $\mathbf{62.4 \pm 2.4\%}$ | $\mathbf{72.4 \pm 5.4\%}$ | $\mathbf{90.2 \pm 4.4\%}$ | $\mathbf{71.2 \pm 4.2\%}$ |
| OpenVLA-OFT (fine-tuned) [18] | − | Clean | $3.8 \pm 0.0\%$ | $1.7 \pm 0.0\%$ | $3.8 \pm 0.0\%$ | $9.3 \pm 0.0\%$ | $4.7 \pm 0.0\%$ |
| | $\Theta$ | Random | $35.6 \pm 3.2\%$ | $38.4 \pm 4.6\%$ | $31.8 \pm 5.8\%$ | $62.2 \pm 10.4\%$ | $42.0 \pm 6.0\%$ |
| | | Best | $\mathbf{56.2 \pm 2.0\%}$ | $\mathbf{64.0 \pm 4.2\%}$ | $\mathbf{67.4 \pm 2.6\%}$ | $\mathbf{82.4 \pm 3.8\%}$ | $\mathbf{67.6 \pm 3.2\%}$ |
| | $\Lambda$ | Random | $16.8 \pm 3.4\%$ | $21.2 \pm 2.8\%$ | $18.0 \pm 4.2\%$ | $46.6 \pm 11.0\%$ | $25.6 \pm 5.4\%$ |
| | | Best | $\mathbf{40.4 \pm 3.0\%}$ | $\mathbf{34.8 \pm 2.4\%}$ | $\mathbf{44.2 \pm 3.6\%}$ | $\mathbf{73.6 \pm 2.8\%}$ | $\mathbf{48.2 \pm 3.0\%}$ |
| | $\phi$ | Random | $11.2 \pm 4.0\%$ | $12.6 \pm 3.6\%$ | $22.4 \pm 5.2\%$ | $38.8 \pm 7.4\%$ | $21.2 \pm 5.0\%$ |
| | | Best | $\mathbf{45.0 \pm 3.2\%}$ | $\mathbf{47.6 \pm 2.2\%}$ | $\mathbf{57.8 \pm 4.4\%}$ | $\mathbf{75.8 \pm 3.4\%}$ | $\mathbf{56.6 \pm 3.3\%}$ |

boundary setting ensures that patches appear prominently in the camera's field of view while remaining within the robot's primary workspace. The $1/3$ margin from texture boundaries prevents patches from being partially cut off or appearing at extreme viewing angles where their impact might be diminished. This placement strategy maximizes the adversarial influence on input images while maintaining realistic scenarios where foreign objects or markings could naturally appear on the work surface.

**Baselines.** To comprehensively evaluate the effectiveness of our methods, we establish three baseline conditions for comparison. First, we assess models on the clean environment, which serves as the control condition where models operate under normal circumstances without any adversarial modifications. Second, we implement a random strategy that samples transformation parameters uniformly from the initial parameter spaces without optimization, quantifying the impact of unguided perturbations and demonstrating the necessity of our optimization-based approach over naive random exploration. Third, we evaluate the best-case attacks using optimized parameters obtained through our proposed optimization framework, representing the most effective parameters discovered for each physical transformation type.

*B. Experiment Setup*

**Dataset & Threat Models.** We conduct our experiments exclusively on the LIBERO dataset [30] in simulation environments, which is a comprehensive simulation dataset designed to evaluate vision-language-action models across four distinct task categories: Spatial, Object, Goal, and Long-horizon tasks. Our method involves physical transformations that require precise control and reproducibility of environmental conditions. Simulation environments enable us to systematically apply and evaluate these physical transformations under controlled conditions, ensuring that the attacks are repeatable and measurable. We select publicly available and state-of-the-art VLAs as victim models for comprehensive evaluation. Specifically, we evaluate eight OpenVLA-based models: four baseline OpenVLA [4] variants fine-tuned on distinct LIBERO task suites (Spatial, Object, Goal, and Long) using standard training procedures, and four corresponding OpenVLA-OFT [18] variants.

**Evaluation Metric.** For task execution assessment, we employ the maximum step count from each LIBERO task suite's training data as the termination criterion for timeout failures, thereby minimizing computational costs. We utilize the Failure Rate (FR), calculated as $1 -$ Success Rate (SR), as established in LIBERO [30], as our principal performance indicator. We also provide the standard deviation of FR across individual tasks within respective task suites to capture performance variability.

**Physical Experiment setting.** We adopt a robotic platform consisting of an AgileX Piper arm equipped with a 1-DOF gripper to provide a 7DoF motion. The sensing system includes a RealSense D435if camera in the fixed position to capture third-view images.

**Evaluation Details.** To evaluate the effectiveness of our proposed methods, we conduct experiments on the LIBERO dataset [30] following the evaluation protocol established by Kim et al. [4]. Each task suite comprises 10 distinct tasks, with 50 trials performed per task, yielding 500 total rollouts per suite. To balance computational efficiency with attack effectiveness, we configure our optimization process with 10 samples per iteration across 50 iterations, enabling thorough exploration of the adversarial parameter space while maintaining reasonable computational overhead. All experiments are performed on a single NVIDIA A800 GPU with 80GB memory.

*C. Main Results*

**Quantitative Results.** Table I presents the failure rates of OpenVLA and OpenVLA-OFT models under various physical transformations across four LIBERO task suites, revealing critical vulnerabilities in VLA. Both models demonstrate substantial performance degradation under all three types, with OpenVLA's average failure rate increasing dramati-
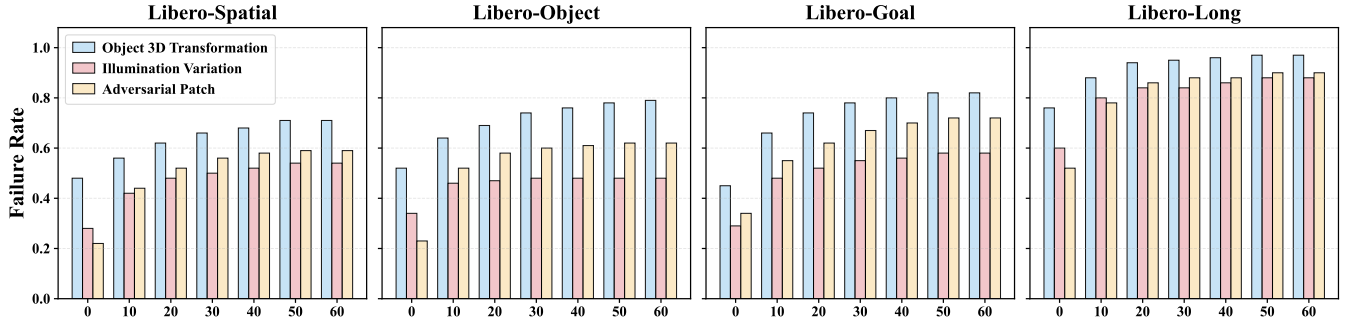
Fig. 4: **Ablation study on the impact of different physical variations across LIBERO benchmark tasks** We evaluate task failure rates under three adversarial conditions: Object 3D Transformation, Illumination Variation, and Adversarial Patch across four LIBERO task suites (Spatial, Object, Goal, and Long). The x-axis represents optimization iterations (0-60), and the y-axis shows the failure rate.

cally from 23.5% (clean) to 82.6% under optimized object transformations ($\Theta$), 62.6% under illumination variations ($\Lambda$), and 71.2% under adversarial patches ($\phi$). Despite OpenVLA-OFT's superior baseline performance with only 4.7% clean failure rate, it remains highly vulnerable to physical attacks, reaching failure rates of 67.6%, 48.2%, and 56.4% for $\Theta$, $\Lambda$, and $\phi$ attacks, respectively. Among these types, object 3D transformations prove to be the most effective, suggesting that spatial reasoning in VLAs is particularly vulnerable to geometric perturbations, while illumination variations and adversarial patches exhibit moderate effectiveness. Notably, even random perturbations without optimization cause significant performance drops (33.2%-55.7% for OpenVLA, 21.2%-42.0% for OpenVLA-OFT), indicating inherent brittleness to environmental variations. This vulnerability is especially pronounced in long-horizon tasks, where failure rates reach 97.8% for OpenVLA and 82.4% for OpenVLA-OFT under optimized object transformations, suggesting that adversarial effects compound over extended action sequences, while simpler Object and Spatial tasks show relatively lower but still substantial vulnerability. These consistent patterns across both standard and optimized models demonstrate that current VLA architectures fundamentally lack robust mechanisms to handle physical variations, particularly in multi-step scenarios.

**Qualitative Results.** We qualitatively analyze robot movement trajectories under the three proposed physical variation types in Fig. 3. For object 3D transformations, as shown in the trajectory plots, the robot maintains similar initial movement patterns but fails to achieve correct object placement due to spatial misalignment. The trajectory demonstrates that while the robot attempts to complete the full task sequence, the geometric perturbations cause mislocalization of target positions, resulting in the robot placing objects at incorrect locations. We attribute this to object 3D transformations of the model's spatial reasoning capabilities through geometric transformations, causing systematic errors in perceiving and reaching intended placement positions. For illumination variations, we observe degraded object recognition that manifests as incomplete grasping actions and premature trajectory termination. The adversarial lighting conditions interfere with

the model's visual perception, leading to failed object detection and misaligned manipulation attempts. For adversarial patches, we observe oscillatory behaviors and intermittent loss of object contact. The patches induce consistent directional biases in the robot's movements, causing it to deviate from intended pushing trajectories and lose task-relevant object interactions. In summary, our qualitative analysis demonstrates that all three physical variations can effectively disrupt robot actions through distinct failure modes.
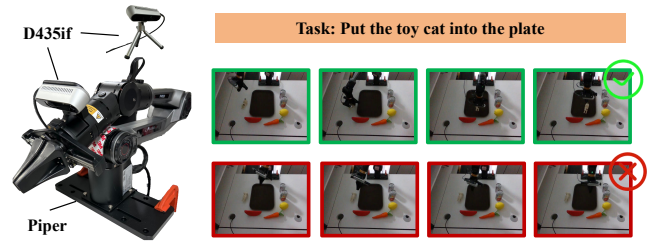


Fig. 5: **Physical Experiment.** The experimental platform consists of an AgileX Piper arm equipped with a 1-DOF gripper and utilize a RealSense D435if camera to capture third-view images. **(Left)** Qualitative Results of the physical world. The first and second row show clean and object 3d transformations case respectively. **(Right)**

**Real-World Performance.** Beyond simulation results, we conduct comprehensive evaluations of the three physical transformations in real-world scenarios using OpenVLA trained on physical world data from BridgeData v2 [31]. The evaluation encompasses three manipulation tasks (object grasping, positioning, and placement), with 50 trials conducted for each physical transformation type, achieving 44.6% average attack success rate. As demonstrated in Fig. 5, when applying 3D transformations to critical objects in the scene, the robot fails to complete designated tasks. Moreover, the unstable movements, manifesting as jerky and oscillatory motions similar to those observed in simulation, pose significant risks to human safety and the operational environment. Additional results are provided in the Appendix.

### D. Ablation Study

We conduct comprehensive ablation experiments analyzing the performance of three physical variations across

different iteration steps on four LIBERO task suites in Fig. 4. All attack types exhibit similar convergence patterns: failure rates increase sharply during the initial iterations (0-20), followed by gradual stabilization, and eventually plateau around iteration 40-50, with Object 3D Transformation exhibiting the steepest degradation. These findings highlight the critical need for developing robust policies, particularly for long-horizon manipulation tasks where cascading errors amplify the impact of adversarial perturbations.

*E. Discussion*

**Defense Mechanisms and Robustness Improvements.** Our findings reveal critical vulnerabilities in current VLA architectures, necessitating comprehensive defense strategies. We recommend incorporating spatial augmentation with systematic rotation and translation variations during training, coupled with multi-view consistency to maintain stable object representations under viewpoint transformation. Illumination-invariant visual encoders pre-trained on diverse illumination conditions can reduce sensitivity to Illumination variations, while training with varied environmental contexts can mitigate adversarial patch effectiveness. Most importantly, the Eva-VLA framework itself provides a direct path to robustness through adversarial training—by using Eva-VLA as an attack generator during training, models can learn to resist worst-case physical variations discovered automatically across simulation environments, effectively transforming our evaluation tool into a defense mechanism that explicitly minimizes sensitivity to physical variations.

## V. CONCLUSIONS

In this paper, we propose Eva-VLA, a novel framework that systematically investigates the vulnerabilities of VLA models, uncovering significant performance degradation under physical variations. Extensive experiments demonstrate that even state-of-the-art VLA models are highly susceptible to these challenges, with issues persisting in both simulated and real-world settings. The observed instability and task failures in physical experiments highlight potential risks to operational safety, underscoring the urgent need for enhanced defenses against physical perturbations. These findings provide critical insights for improving the robustness of VLA-based robotic systems in adversarial environments.

## REFERENCES

[1] Y. Ma, Z. Song, Y. Zhuang, J. Hao, and I. King, "A survey on vision-language-action models for embodied ai," *arXiv preprint arXiv:2405.14093*, 2024.

[2] Q. Li, Y. Liang, Z. Wang, L. Luo, X. Chen, M. Liao, F. Wei, Y. Deng, S. Xu, Y. Zhang, *et al.*, "Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation," *arXiv preprint arXiv:2411.19650*, 2024.

[3] S. Li, J. Wang, R. Dai, W. Ma, W. Y. Ng, Y. Hu, and Z. Li, "Robonurse-vla: Robotic scrub nurse system based on vision-language-action model," *arXiv preprint arXiv:2409.19590*, 2024.

[4] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, *et al.*, "Open-vla: An open-source vision-language-action model," *arXiv preprint arXiv:2406.09246*, 2024.

[5] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, *et al.*, "π0: A vision-language-action flow model for general robot control. corr, abs/2410.24164, 2024. doi: 10.48550," *arXiv preprint ARXIV.2410.24164*.

[6] T. Wang, D. Liu, J. C. Liang, W. Yang, Q. Wang, C. Han, J. Luo, and R. Tang, "Exploring the adversarial vulnerabilities of vision-language-action models in robotics," *arXiv preprint arXiv:2411.13587*, 2024.

[7] N. Hansen, "The cma evolution strategy: A tutorial," *arXiv preprint arXiv:1604.00772*, 2016.

[8] X. Li, M. Liu, H. Zhang, C. Yu, J. Xu, H. Wu, C. Cheang, Y. Jing, W. Zhang, H. Liu, *et al.*, "Vision-language foundation models as effective robot imitators," *arXiv preprint arXiv:2311.01378*, 2023.

[9] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, *et al.*, "Rt-2: Vision-language-action models transfer web knowledge to robotic control," *arXiv preprint arXiv:2307.15818*, 2023.

[10] R. Zheng, Y. Liang, S. Huang, J. Gao, H. Daumé III, A. Kolobov, F. Huang, and J. Yang, "Tracevla: Visual trace prompting enhances spatial-temporal awareness for generalist robotic policies," *arXiv preprint arXiv:2412.10345*, 2024.

[11] X. Li, P. Li, M. Liu, D. Wang, J. Liu, B. Kang, X. Ma, T. Kong, H. Zhang, and H. Liu, "Towards generalist robot policies: What matters in building vision-language-action models," *arXiv preprint arXiv:2412.14058*, 2024.

[12] J. Zheng, J. Li, D. Liu, Y. Zheng, Z. Wang, Z. Ou, Y. Liu, J. Liu, Y.-Q. Zhang, and X. Zhan, "Universal actions for enhanced embodied foundation models," *arXiv preprint arXiv:2501.10105*, 2025.

[13] K. Pertsch, K. Stachowicz, B. Ichter, D. Driess, S. Nair, Q. Vuong, O. Mees, C. Finn, and S. Levine, "Fast: Efficient action tokenization for vision-language-action models," *arXiv preprint arXiv:2501.09747*, 2025.

[14] X. Chen, J. Djolonga, P. Padlewski, B. Mustafa, S. Changpinyo, J. Wu, C. R. Ruiz, S. Goodman, X. Wang, Y. Tay, *et al.*, "Pali-x: On scaling up a multilingual vision and language model," *arXiv preprint arXiv:2305.18565*, 2023.

[15] S. Karamcheti, S. Nair, A. Balakrishna, P. Liang, T. Kollar, and D. Sadigh, "Prismatic vlms: Investigating the design space of visually-conditioned language models," in *Forty-first International Conference on Machine Learning*, 2024.

[16] Q. Vuong, S. Levine, H. R. Walke, K. Pertsch, A. Singh, R. Doshi, C. Xu, J. Luo, L. Tan, D. Shah, *et al.*, "Open x-embodiment: Robotic learning datasets and rt-x models," in *Towards Generalist Robots: Learning Paradigms for Scalable Skill Acquisition@ CoRL2023*, 2023.

[17] D. Qu, H. Song, Q. Chen, Y. Yao, X. Ye, Y. Ding, Z. Wang, J. Gu, B. Zhao, D. Wang, *et al.*, "Spatialvla: Exploring spatial representations for visual-language-action model," *arXiv preprint arXiv:2501.15830*, 2025.

[18] M. J. Kim, C. Finn, and P. Liang, "Fine-tuning vision-language-action models: Optimizing speed and success," *arXiv preprint arXiv:2502.19645*, 2025.

[19] A. Du, B. Chen, T.-J. Chin, Y. W. Law, M. Sasdelli, R. Rajasegaran, and D. Campbell, "Physical adversarial attacks on an aerial imagery object detector," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 1796–1806.

[20] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.

[21] A. Liu, X. Liu, J. Fan, Y. Ma, A. Zhang, H. Xie, and D. Tao, "Perceptual-sensitive gan for generating adversarial patches," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 1028–1035.

[22] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.

[23] J. Su, D. V. Vargas, and K. Sakurai, "One pixel attack for fooling deep neural networks," *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 5, pp. 828–841, 2019.

[24] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *International conference on machine learning*. PMLR, 2018, pp. 284–293.

[25] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," *arXiv preprint arXiv:1712.09665*, 2017.

[26] H. Liu, S. Ruan, Y. Huang, S. Zhao, and X. Wei, "When lighting deceives: Exposing vision-language models' illumination vulnerability through illumination transformation attack," *arXiv preprint arXiv:2503.06903*, 2025.

[27] S. Ruan, H. Liu, Y. Huang, X. Wang, C. Kang, H. Su, Y. Dong, and X. Wei, "Advdreamer unveils: Are vision-language models truly ready for real-world 3d variations?" *arXiv preprint arXiv:2412.03002*, 2024.

[28] S. Haddadin, S. Parusel, L. Johannsmeier, S. Golz, S. Gabl, F. Walch, M. Sabaghian, C. Jähne, L. Hausperger, and S. Haddadin, "The franka emika robot: A reference platform for robotics research and education," *IEEE Robotics & Automation Magazine*, vol. 29, no. 2, pp. 46–64, 2022.

[29] D. Golovin, B. Solnik, S. Moitra, G. Kochanski, J. Karro, and D. Sculley, "Google vizier: A service for black-box optimization," in *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, 2017, pp. 1487–1495.

[30] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone, "Libero: Benchmarking knowledge transfer for lifelong robot learning," *Advances in Neural Information Processing Systems*, vol. 36, pp. 44776–44791, 2023.

[31] H. R. Walke, K. Black, T. Z. Zhao, Q. Vuong, C. Zheng, P. Hansen-Estruch, A. W. He, V. Myers, M. J. Kim, M. Du, *et al.*, "Bridgedata v2: A dataset for robot learning at scale," in *Conference on Robot Learning*. PMLR, 2023, pp. 1723–1736.