

# PACT: Self-Evolving Physical Safety Alignment for Diffusion Policies in Embodied Manipulation

Anonymous Authors<sup>1</sup>

## Abstract

Diffusion policies have achieved remarkable success in robotic manipulation, yet they often fail to satisfy strict physical constraints required for safe deployment. Existing approaches impose safety either prematurely during training or reactively via external guardrails at test time, limiting policy expressivity and overall scalability. We propose Physical safety Alignment for Constrained Trajectories (PACT), a self-evolving post-training framework that projects pretrained diffusion policies onto constraint-feasible regions without accessing demonstration data or task rewards. PACT distills constraint gradients into the diffusion model through a reverse-KL objective with dense supervision across timesteps. It incorporates a curriculum that progressively tightens constraints while maintaining theoretically bounded policy shift and monotone improvement, mitigating the safety-performance trade-off from catastrophic forgetting. On simulated and real-world embodied manipulation benchmarks, PACT significantly reduces safety violations by 31.0% on average while improving task success by 30.7%.

## 1. Introduction

Recent advances in diffusion policies have demonstrated remarkable progress in robotic manipulation (Janner et al., 2022; Pearce et al., 2023; Chi et al., 2025), enabling expressive, multimodal action distributions and strong generalization across diverse tasks and environments (Liu et al., 2025b; Amin et al., 2025). Despite their power, real-world deployment requires strict adherence to physical constraints, such as collision avoidance, force limits, and kinematic feasibility (Rueß, 2022). Constraint violations can be catastrophic, leading to task failure, irreversible hardware damage, and

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

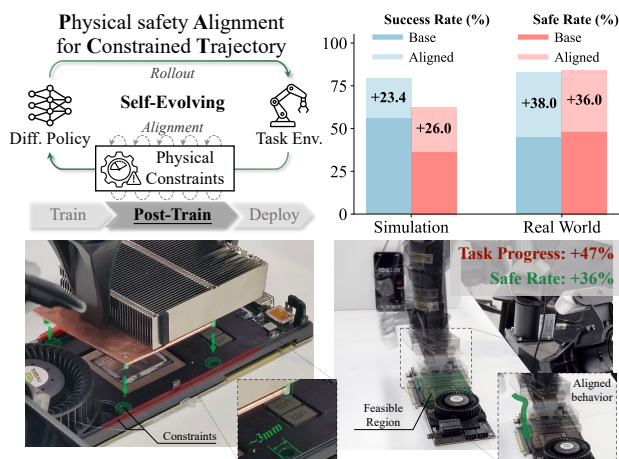


Figure 1. **Physical safety alignment for diffusion-based manipulation.** PACT aligns a pretrained diffusion policy in a post-training stage using self-rollouts and continuous constraint supervision throughout the diffusion process in a *self-evolving manner without external demonstrations or rewards*. *Top right*: PACT improves both task performance and safety in simulation and real-world settings, resolving the safety–performance trade-off. *Bottom*: GPU assembly requiring millimeter-level precision, PACT enables emergent fine-grained correction behaviors that adjust trajectories into feasible regions, ensuring accurate and safe insertions.

human injury, especially in safety-critical applications like manufacturing (Buhl et al., 2019) and surgery (Hu et al., 2023). Consequently, deploying diffusion policies at scale necessitates preserving their expressive power while guaranteeing constraint satisfaction under all operating conditions.

To mitigate this issue, existing methods can be broadly categorized into design-time (or train-time) and test-time approaches. Design-time approaches (Ross et al., 2011) embed safety constraints during data collection (Liu et al., 2024; Bahety et al., 2025) or model training (García & Fernández, 2015; Ciftci et al., 2025), ensuring robustness to perturbations during execution. However, they require extensive human expertise and task-specific engineering, limiting scalable deployment. Test-time approaches (Hsu et al., 2023) instead employ guardrails (Reichlin et al., 2022; Deng et al., 2025) or backup policies to intercept unsafe actions during execution (Wong et al., 2022). However, they depend on predefined constraints or additional sensors for real-time monitoring. Therefore, current methods impose safety constraints either prematurely, restricting training and

potentially impeding learning, or reactively through limited deployment-stage interventions. Both strategies compromise policy expressivity and scalability while failing to provide principled safety guarantees for deployment.

Drawing inspiration from post-training safety alignment in large language models (Bai et al., 2022; Dai et al., 2024; Chen et al., 2024), we align pretrained diffusion policies for physical safety before deployment. We formulate it as projecting pretrained policies onto the feasible physically constrained regions in a Constrained Markov Decision Process (CMDP) (Altman, 2021). This preserves the expressiveness of the original policies, enables flexible adaptation to varying constraint sets without retraining, and maintains computational efficiency by decoupling safety from capability learning. However, safety alignment for embodied manipulation (Billard & Kragic, 2019) presents unique challenges. First, diffusion models inherently spread probability mass across the action space, complicating strict constraint satisfaction even after alignment. Second, post-training cannot access demonstration data or task utility functions, precluding conventional fine-tuning through data recollection (Liu et al., 2024) or reward-based optimization (Garcia & Fernández, 2015). Third, aggressive safety enforcement risks inducing catastrophic forgetting (McCloskey & Cohen, 1989), since restricting policy support may eliminate task-critical behaviors when safe and original distributions diverge. Collectively, these challenges require scalable post-training mechanisms that enforce physical constraints while preserving the expressivity and the task performance.

To address these challenges, we propose **Physical safety Alignment for Constrained Trajectories (PACT)**, a self-evolving post-training framework that aligns diffusion policies with physical safety constraints. First, PACT introduces constraint distillation to tackle mode concentration by densely injecting gradient signals from safety cost functions at every diffusion timestep. This is equivalent to minimizing the backward KL divergence that concentrates the policy distribution into feasible regions while mitigating exposure bias (Bengio et al., 2015). In contrast to RL methods that suffer from sparse or noisy rewards, PACT provides stable, continuous supervision throughout the entire diffusion process. Second, self-evolving optimization overcomes the data accessibility challenge by operating exclusively on self-generated trajectories, eliminating dependencies on external demonstrations or reward engineering. Third, we incorporate a curriculum-based distillation that progressively tightens safety constraints to mitigate the performance-safety trade-off (Zhang et al., 2025b). This curriculum theoretically promotes bounded policy shift and monotonic improvement, preventing abrupt action collapse and preserving task competence by avoiding catastrophic forgetting. Together, these components constitute a scalable post-training alignment framework that embeds physical safety into diffusion

policies through curriculum-based constraint distillation on self-collected rollouts, substantially reducing human effort by eliminating the need for demonstrations, task-specific utility functions, or outcome annotations.

We evaluate the policy aligned using PACT across simulated and real-world manipulation benchmarks that require delicate bimanual coordination, *without privileged state information or explicit constraint computation* at test time. PACT reduces safety violation rates by 31.0% while improving task success rates by 30.7% and achieving higher training efficiency compared to RL-based safety alignment methods. Notably, PACT exhibits strong real-world performance, including challenging GPU assembly tasks requiring millimeter-level precision, validating PACT as a practical paradigm for deploying diffusion policies in safety-critical robotic systems. Our contributions are as follows:

- To our knowledge, we introduce the first post-training framework for aligning diffusion-based manipulation policies with physical safety constraints, requiring no manual data collection or test-time guardrails.
- We develop a constraint distillation framework with curriculum-based updates that preserves task performance during safety enforcement in the absence of task utilities or demonstrations, enabling stable on-policy optimization with provable, continuous safety improvement.
- Our method substantially reduces safety violations while improving task success rates and training efficiency relative to RL-based alignment, and demonstrates real-world performance on challenging manipulation tasks.

## 2. Background

**Diffusion and flow policies.** Explicitly modeling stochasticity and multi-modality, diffusion-based policies (Sohl-Dickstein et al., 2015; Ho et al., 2020) enable more expressive control representations and have proven effective for real-world robotic tasks involving diverse behaviors (Janner et al., 2022; Chen et al., 2023; Chi et al., 2025). Given an action  $\mathbf{a}$ , a forward corruption process gradually injects Gaussian noise with a specific noise schedule  $\alpha_t$  and  $\sigma_t$ :

$$\mathbf{a}_t = \alpha_t \mathbf{a} + \sigma_t \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), \quad t \in [0, 1]. \quad (1)$$

Intuitively, diffusion policies are learned by training a denoising model  $\epsilon_\phi$  that predicts the injected noise  $\epsilon$ :

$$\min_{\phi} \mathbb{E}_{t, \epsilon, \mathbf{s}, \mathbf{a} \sim \mu(\cdot | \mathbf{s})} [\|\epsilon_\phi(\mathbf{a}_t, \mathbf{s}, t) - \epsilon\|_2^2], \quad (2)$$

which implicitly recovers the score function of the diffused behavior distribution (Song et al., 2021c):

$$\nabla_{\mathbf{a}_t} \mu_t(\mathbf{a}_t | \mathbf{s}, t) \approx \nabla_{\mathbf{a}_t} \mu_\phi(\mathbf{a}_t | \mathbf{s}, t) = \epsilon_\phi(\mathbf{a}_t, \mathbf{s}, t) / \sigma_t, \quad (3)$$

and allows sampling actions from the learned behavior policy  $\mu_\phi$ . Beyond discrete denoising formulations, flow policies parameterize the policy as a velocity field that defines a deterministic flow transporting samples from noise to the action distribution (Liu et al., 2023), proven to be an equivalent continuous-time interpretation through probability flow ODEs (Lipman et al., 2023; Song et al., 2021a). Below we view diffusion and flow-based behavior policies as a unified generative policy class (Kingma & Gao, 2023) for subsequent physical safety alignment.

**Physical safety enforcement.** Prior work on physical safety in robotic manipulation (Haddadin, 2015) can be categorized by the stage at which safety is enforced. Development-time approaches incorporate safety during data collection or policy learning (Brunke et al., 2022), primarily through imitation learning (IL) (Ross et al., 2011; Yang et al., 2024; Liu et al., 2024; Bahety et al., 2025) and safe reinforcement learning (Safe RL) (Garcia & Fernández, 2015; Altman, 2021; Ying et al., 2022). IL-based methods imitate curated safe demonstrations but suffer from degraded performance due to restricted action support. Safe RL formulates safety as auxiliary cost functions or constrained objectives, and typically relies on task utility to drive learning (Altman, 2021; Ying et al., 2022; Lee et al., 2022). However, the rewards are sparse and delayed in real-world settings (Gu et al., 2024), resulting in poor performance. Inference-time methods enforce safety through action filtering or control barrier functions that intervene upon constraint violations (Reichlin et al., 2022; Wabersich et al., 2023; Deng et al., 2025). However, these approaches require privileged state access, additional sensors, or perception models (Wong et al., 2022; Gokmen et al., 2023), thereby increasing inference-time overhead and limiting applicability to novel scenarios. Post-training methods, recently explored by SafeVLA (Zhang et al., 2025a), amortize safety supervision into the policy at the intermediate stage between capability acquisition and deployment. However, it is limited to discrete categorical distribution and low-dimensional navigation tasks. We address *diffusion policies* operating over *continuous, high-dimensional* action spaces with complex physical safety constraints for *manipulation*, posing substantial challenges.

### 3. Methodology

#### 3.1. Problem Formulation

We study *physical safety alignment* of diffusion policies under explicit safety constraints. The environment is modeled as a *constrained Markov Decision Process* (CMDP)  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{C} \rangle$  (Altman, 2021), where  $\mathcal{S}$  and  $\mathcal{A}$  respectively denote the state and action spaces,  $\mathcal{P}(s_{t+1} | s_t, \mathbf{a}_t)$  denotes the transition dynamics,  $\mathcal{C} = \{c_k(s, \mathbf{a})\}_{k=1}^m$  is a set of safety cost functions encoding physical constraints (e.g., collisions or excessive force). The *safe region* (feasible

policy set) derived from safety costs  $\mathcal{C}$  is defined as

$$\Pi_{\text{safe}} = \left\{ \pi : \mathbb{E}_{(s, \mathbf{a}) \sim d^\pi} c_k(s, \mathbf{a}) \leq d_k, \forall k \right\}. \quad (4)$$

Unlike standard CMDP formulations that rely on rewards, we frame *safety alignment* as a *regularized projection* of the base policy  $\mu_\phi(\mathbf{a}_t | s_t)$  onto the feasible set  $\Pi_{\text{safe}}$ :

$$\arg \min_{\pi \in \Pi_{\text{safe}}} \mathbb{E}_{s \sim d^\pi} [D_{\text{KL}}(\pi(\cdot | s) \| \mu_\phi(\cdot | s))], \quad (5)$$

which seeks the closest safety-compliant policy while preserving the behavioral fidelity of the original policy. Furthermore, we assume *no access to the pretrained behavior dataset* during safety alignment, which is realistic but restrictive (Dulac-Arnold et al., 2021). This precludes replay-based regularization or supervised correction to control distribution shift, making direct constrained optimization infeasible (Ball et al., 2023). To solve Problem (5), we leverage the Lagrange multiplier method to incorporate safety constraints while minimizing the KL divergence from the base policy (Chow et al., 2018; Tessler et al., 2018). Specifically, we introduce a set of Lagrange multipliers  $\lambda_k \geq 0$  for each safety constraint  $c_k(s, \mathbf{a})$  which enforce the safety conditions in the optimization process. The Lagrangian  $\mathcal{L}(\pi, \lambda)$  for the constrained optimization problem is formulated as:

$$\mathbb{E}_{(s, \mathbf{a}) \sim d^\pi} [D_{\text{KL}}(\pi(\cdot | s) \| \mu_\phi(\cdot | s)) + \sum_{k=1}^m \lambda_k (c_k(s, \mathbf{a}) - d_k)], \quad (6)$$

which shares a similar structure with objective in Regularized RL (Dayan & Hinton, 1997; Schulman et al., 2017).

#### 3.2. Physical Safety Alignment for Diffusion Policies

To minimize the objective in Eq. (6), extensive prior work has studied parametric policies such as Gaussian or categorical distributions using policy gradient (Schulman et al., 2017). However, extending these approaches to diffusion policies is non-trivial. In particular, policy gradient methods require tractable likelihood ratios or entropy terms, whereas diffusion policies only admit implicit densities whose likelihoods can only be computed approximately via costly probability ODE solvers or variational bounds for SDEs (Song et al., 2021b). These limitations make direct optimization of Eq. (6) impractical for diffusion models. Despite the intractability of likelihood-based optimization, the optimal solution to Eq. (6) admits a simple closed-form structure. For fixed Lagrange multipliers  $\{\lambda_k\}_{k=1}^m$ , the constrained optimal policy satisfies (Peng et al., 2019):

$$\pi^*(\mathbf{a} | s) \propto \mu_\phi(\mathbf{a} | s) \exp\left(-\sum_{k=1}^m \lambda_k c_k(s, \mathbf{a})\right), \quad (7)$$

which can be interpreted as an exponential tilting of the base policy by safety costs. Although the normalized density in

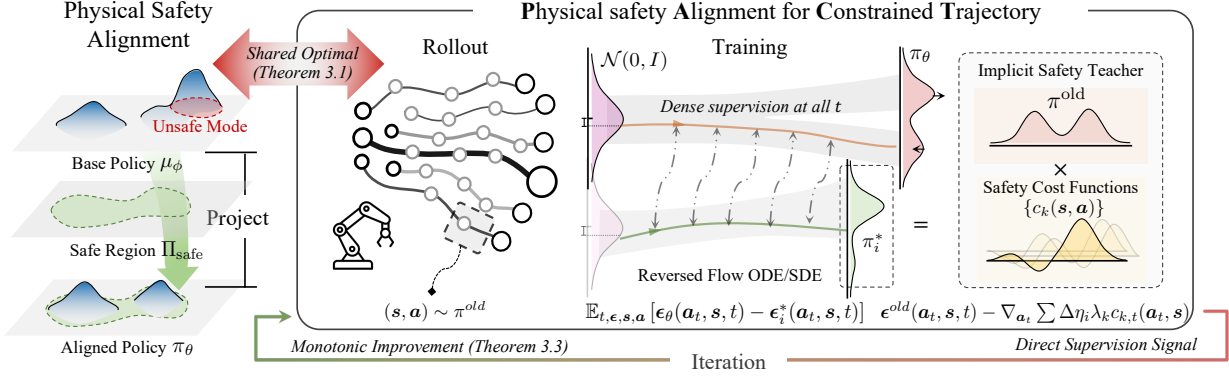


Figure 2. **Overview of Physical safety Alignment for Constrained Trajectories.** PACT frames physical safety alignment as projecting a pretrained diffusion policy  $\mu_\phi$  onto the CMDP feasible set  $\Pi_{\text{safe}}$  via a KL-regularized constrained objective. For fixed multipliers, the score function of the optimal aligned policy is defined as an *implicit safety teacher* by combining the base score with differentiable cost gradients (Theorem 3.1). We distill this teacher into a student diffusion policy  $\pi_\theta$  using self-rollouts with direct and dense supervision across diffusion time and in a self-evolving manner without manual demonstrations or rewards. A curriculum schedule progressively increases constraint strength, yielding controlled policy shifts and monotonic safety improvement over iterations (Theorem 3.3).

Eq. (7) remains intractable, it is available for sampling as its score function is directly computable:

$$\epsilon^*(\mathbf{a}_t, \mathbf{s}, t) \triangleq \epsilon_\phi(\mathbf{a}_t, \mathbf{s}, t) - \sum_{k=1}^m \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t) / \sigma_t, \quad (8)$$

where  $c_{k,t}(\cdot)$  denotes the intermediate costs derived from the original cost function  $c_k(\cdot)$  (Lu et al., 2023) with its detailed structure elaborated in Appx. A.1. We define the score function  $\epsilon^*(\cdot)$  in Eq. (8) as the *implicit safety teacher*.

However, sampling with  $\epsilon^*(\cdot)$  requires evaluating the cost functions  $c_k$ , which often depend on privileged environment states (e.g., object poses) typically obtained via additional sensors or auxiliary perception modules. This reliance hinders direct test-time deployment and is inherently unscalable, as the per-agent sensing and computation overhead grows with the size of the robot fleet. To overcome this limitation, we approximate  $\epsilon^*$  via distillation by training a student diffusion policy  $\epsilon_\theta$  to match the implicit teacher’s score function for physical safety alignment<sup>1</sup>:

$$\min_{\theta} \mathbb{E}_{t, \epsilon, (\mathbf{s}, \mathbf{a}) \sim d^{\pi_\theta}} \|\epsilon_\theta(\mathbf{a}_t, \mathbf{s}, t) - \epsilon^*(\mathbf{a}_t, \mathbf{s}, t)\|^2. \quad (9)$$

**Theorem 3.1** (Optimality of PACT). *Given unlimited model capacity and sufficient data, the optimal solution for the distillation objective in Eq. (9) is the score function of Eq. (7)*<sup>2</sup>.

**Remark 3.2.** This theorem shows that distillation can recover the optimal solution of the constrained optimization problem without explicit likelihood evaluation. Compared to RL-based methods that rely on stochastic, high-variance reward signals, it provides direct and stable guidance toward the feasible region without requiring outcome rewards or

<sup>1</sup>The objective in Eq. (9) is parameterization-agnostic. Further discussion is provided in Appx. A.2.

<sup>2</sup>Proof in Appx. A.3

value models. Moreover, the proposed objective is *solver-agnostic*, enabling efficient sampling with few denoising steps and higher-order solvers during rollout, unlike prior diffusion RL methods that are tightly coupled to first-order SDE samplers (Black et al., 2024b; Liu et al., 2025a).

**Curriculum distillation to mitigate irreversible Out-of-Distribution collapse.** In practice, a key failure mode occurs when transient intermediate policies deviate substantially from the base behavior, inducing rapid rollout state distribution drift and triggering *Irreversible Out-of-Distribution (OOD) Collapse*, wherein an under-optimized policy is driven into OOD regimes and suffers a severe loss of task competence. This issue is exacerbated in safety alignment because task utility is inaccessible; therefore, no corrective supervision is available to steer recovery, making the degradation effectively irreversible. Mechanistically, it is primarily due to the lack of trust region control over intermediate updates (Conn et al., 2000) for direct distillation, so small deviations can compound over time and rapidly push trajectories into OOD states. To address this, we introduce a curriculum distillation strategy that progressively enforces safety constraints, enabling a smooth transition toward the final aligned solution. Concretely, we adopt a monotonically increasing curriculum schedule  $\eta_0, \dots, \eta_N \in [0, 1]$  with  $\eta_0 = 0$  and  $\eta_N = 1$ , which gradually scales the constraint multipliers over  $N$  iterations. At each iteration  $i$ , we solve the following intermediate objective:

$$\min_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim d^{\pi}} \left[ D_{\text{KL}}(\pi(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) + \sum_{k=1}^m \Delta \eta_i \cdot \lambda_k (c_k(\mathbf{s}, \mathbf{a}) - d_k) \right], \quad (10)$$

where  $\Delta \eta_i = \eta_i - \eta_{i-1}$  ( $i > 0$ ) denotes the incremental schedule step.  $\pi^{\text{old}}$  is the trained  $\pi$  of the previous iteration, and  $\pi^{\text{old}} = \mu_\phi$  at the initial iteration ( $i = 1$ ). The

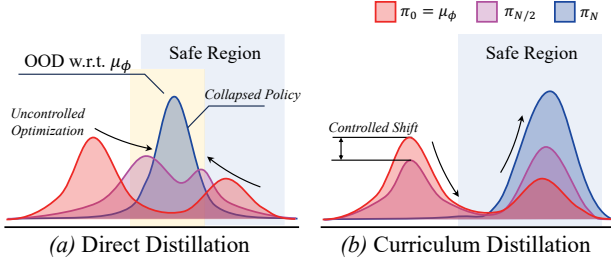


Figure 3. Curriculum distillation mitigates Irreversible OOD Collapse by controlling intermediate policy shift. We illustrate the evolution of policy distributions over iterations: (a) *Direct distillation* enforces constraints without control, so intermediate policies can drift rapidly, pushing rollouts into OOD regions and yielding a collapsed policy that loses task competence despite aiming for safety. (b) *Curriculum distillation* progressively increases constraint strength and regularizes policy within a trust region, inducing controlled transition to a safety-compliant policy.

corresponding distillation objective at iteration  $i$  is:

$$\min_{\theta} \mathbb{E}_{t, \epsilon, (\mathbf{s}, \mathbf{a}) \sim d^{\pi_{\theta}}} \|\epsilon_{\theta}(\mathbf{a}_t, \mathbf{s}, t) - \epsilon_i^+(\mathbf{a}_t, \mathbf{s}, t)\|^2, \text{ with}$$

$$\epsilon_i^+(\mathbf{a}_t, \mathbf{s}, t) = \epsilon^{\text{old}}(\mathbf{a}_t, \mathbf{s}, t) - \sum_{k=1}^m \Delta \eta_i \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t) / \sigma_t, \quad (11)$$

where  $\epsilon_i^+(\cdot)$  defines the *Curriculum Implicit Safety Teacher*, serving as an intermediate target that smoothly interpolates toward the final Implicit Safety Teacher  $\epsilon^*(\cdot)$  in Eq. (8).

**Theorem 3.3** (Monotonic Improvement under Curriculum Distillation). *Assume each cost function is bounded, i.e.,  $c_k(\mathbf{s}, \mathbf{a}) \in [c_k^{\min}, c_k^{\max}]$  for all  $k$  and all  $(\mathbf{s}, \mathbf{a})$ . Define the aggregated Lagrangian cost*

$$\ell(\mathbf{s}, \mathbf{a}) \triangleq \sum_{k=1}^m \lambda_k (c_k(\mathbf{s}, \mathbf{a}) - d_k), \quad \ell(\mathbf{s}, \mathbf{a}) \in [\ell_{\min}, \ell_{\max}].$$

Let  $\pi^{\text{new}}$  be obtained by solving the objective in Eq. (11) at any iteration  $i$ ; the update satisfies:

- **Monotonic improvement.** The expected curriculum Lagrangian cost under  $d^{\pi^i}$  is non-increasing:

$$\mathbb{E}_{d^{\pi^{\text{new}}}}[\ell(\mathbf{s}, \mathbf{a})] \leq \mathbb{E}_{d^{\pi^{\text{old}}}}[\ell(\mathbf{s}, \mathbf{a})]. \quad (12)$$

- **Controlled policy shift.** Set  $\Delta \ell = \ell_{\max} - \ell_{\min}$ , then the policy change is bounded as

$$\mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{KL}}(\pi^{\text{new}}(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) \right] \leq \Delta \eta_i \Delta \ell. \quad (13)$$

Thus, the proposed curriculum distillation ensures monotonic safety improvement along the iterations and constructs a smooth policy with controlled policy shifts to prevent abrupt distributional shifts in the original objective.

**Corollary 3.4** (Curriculum PACT Optimal Solution). *The optimal solution  $\pi_i^*$  at iteration  $i$  satisfies*

$$\pi_i^*(\mathbf{a} | \mathbf{s}) \propto \mu(\mathbf{a} | \mathbf{s}) \exp\left(-\sum_{k=1}^m \eta_i \cdot \lambda_k c_k(\mathbf{s}, \mathbf{a})\right), \quad (14)$$

### Algorithm 1 The Training Pipeline of PACT

**Require:** Pretrained diffusion policy  $\epsilon_{\phi}$ , number of iterations  $N$ , number of epochs  $E$ , set of differentiable safety cost functions  $\{c_k(\mathbf{s}, \mathbf{a})\}_{k=1}^m$  and corresponding fixed Lagrange multipliers  $\{\lambda_k\}_{k=1}^m$ , curriculum schedule  $\{\eta_i\}_{i=0}^N$ , max diffusion time to inject cost guidance  $t_c$ .

**Ensure:** The aligned diffusion policy  $\epsilon_{\theta}$ .

```

1: for iteration  $i \leftarrow 1$  to  $N$  do
2:    $\triangleright$  Rollout and collect data ◁
3:   Rollout policy  $\epsilon_{\theta}$  and collect trajectories as  $\mathcal{D} = \{(\mathbf{s}, \mathbf{a})\}$ .
4:   for epoch  $\leftarrow 1$  to  $E$  do
5:     for all mini-batch  $\{\mathbf{s}, \mathbf{a}\} \in \mathcal{D}$  do
6:        $\mathbf{a}_t = \alpha_t \mathbf{a} + \sigma_t \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ ,  $t \in [0, 1]$ .
7:        $\mathbf{a}_{0|t} = (\mathbf{a}_t - \sigma_t \epsilon^{\text{old}}(\mathbf{s}, \mathbf{a}_t, t)) / \alpha_t$ . (Eq.(1))
8:       Update  $\theta$  by minimizing the objective in Eq. (11)
          with approximated score for Curriculum Implicit
          Safety Teacher computed with Eq. (15).
9:     end for
10:  end for
11:  Update  $\theta_{\text{old}} \leftarrow \theta$  and clear buffer  $\mathcal{D}_{\tau} \leftarrow \emptyset$ 
12: end for
    
```

and the final solution at iteration  $N$  exactly coincides with the shared closed-form optimal solution in Theorem 3.1 given unlimited model capacity and data samples<sup>3</sup>.

In conclusion, the curriculum distillation procedure preserves the optimal solution of the original constrained objective. Its scheduled formulation ensures that each intermediate update enforces only a controlled deviation from the current policy with guaranteed improvement. This gradual enforcement of constraints prevents the policy from being prematurely driven into unfamiliar states, thereby avoiding Irreversible OOD Collapse and ensuring training stability.

### 3.3. Practical Implementations

**Training-free approximation for intermediate cost gradient and few-steps distillation.** The exact gradient of intermediate costs requires an additional denoising-time condition and training steps (Lu et al., 2023). However, injecting cost gradient only during the final few denoising steps of guided sampling is effective in most scenarios (Bao et al., 2022; Xu et al., 2023). Motivated by the few-step nature of guided sampling, we propose a training-free approximation to the intermediate cost gradient based on a Taylor expansion for  $\nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t)$  around  $t = 0$ . Specifically, for sufficiently small  $t$ , we approximate the gradient at  $\mathbf{a}_t$  by evaluating it at the posterior mean of the clean action, defined as  $\mathbf{a}_{0|t} \triangleq \mathbb{E}_{q(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})}[\mathbf{a}_0]$ , yielding  $\nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t) \approx \nabla_{\mathbf{a}_t} c_k(\mathbf{s}, \mathbf{a}_{0|t})$ , where  $q(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})$  denotes the posterior for the clean action  $\mathbf{a}_0$  and  $\mathbf{a}_{0|t}$  can be simply computed by re-formatting Eq. (1) to a one-step reverse diffusion update (Chung et al., 2023). This approximation is widely used in training-free guided sampling

<sup>3</sup>Proof in Appx. A.5

Table 1. Task success and safety across bimanual manipulation benchmarks in RoboTwin. We report *Success Rate (succ.)* and *Safe Rate (safe)* for each base policy before and after post-training with PACT.

Model	Pick Dual Bottle		Pick Diverse Bottle		Handover Apple		Handover Block		Place Dual Shoes		Pour Water		Stack Blocks		Average	
	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe
DP	67%	36%	36%	6%	63%	64%	28%	70%	23%	66%	38%	20%	63%	24%	45%	41%
+ PACT	<b>96%</b>	<b>89%</b>	<b>58%</b>	<b>21%</b>	<b>81%</b>	<b>82%</b>	<b>72%</b>	<b>78%</b>	<b>41%</b>	<b>73%</b>	<b>86%</b>	<b>64%</b>	<b>80%</b>	<b>48%</b>	<b>73%</b>	<b>65%</b>
DP3	44%	24%	4%	2%	54%	46%	59%	71%	9%	59%	16%	13%	6%	6%	27%	32%
+ PACT	<b>60%</b>	<b>46%</b>	<b>11%</b>	<b>7%</b>	<b>70%</b>	<b>59%</b>	<b>72%</b>	<b>73%</b>	<b>19%</b>	<b>70%</b>	<b>41%</b>	<b>16%</b>	<b>27%</b>	<b>16%</b>	<b>43%</b>	<b>41%</b>
RDT-1B	61%	27%	35%	12%	76%	52%	64%	61%	42%	72%	52%	17%	60%	13%	56%	36%
+ PACT	<b>88%</b>	<b>74%</b>	<b>71%</b>	<b>43%</b>	<b>83%</b>	<b>59%</b>	<b>75%</b>	<b>82%</b>	<b>63%</b>	<b>82%</b>	<b>85%</b>	<b>72%</b>	<b>89%</b>	<b>24%</b>	<b>79%</b>	<b>62%</b>
$\pi_{0.5}$	64%	20%	37%	8%	71%	55%	66%	64%	41%	70%	48%	14%	67%	25%	56%	37%
+ PACT	<b>89%</b>	<b>72%</b>	<b>76%</b>	<b>42%</b>	<b>86%</b>	<b>61%</b>	<b>74%</b>	<b>79%</b>	<b>70%</b>	<b>79%</b>	<b>79%</b>	<b>72%</b>	<b>92%</b>	<b>30%</b>	<b>81%</b>	<b>62%</b>

across diverse tasks (Chung et al., 2023; Yu et al., 2023), including bimanual manipulation (Deng et al., 2025). The training-free approximation yields the few-step Curriculum Implicit Safety Teacher in Eq. (11) as:

$$\epsilon_i^+(\mathbf{a}_t, \mathbf{s}, t) \approx \begin{cases} \epsilon^{\text{old}}(\mathbf{a}_t, \mathbf{s}, t) - \sum_{k=1}^m \Delta \eta_i \lambda_k \nabla_{\mathbf{a}_t} c_k(\mathbf{s}, \mathbf{a}_{0:t}), & t < t_c, \\ \epsilon^{\text{old}}(\mathbf{a}_t, \mathbf{s}, t), & t \geq t_c. \end{cases} \quad (15)$$

This requires only the original differentiable cost functions. By restricting constraint injection to the final few steps (e.g.,  $t_c = 0.03$  in all experiments), the approach substantially reduces the computational overhead from the Jacobian evaluations. The theoretical and empirical justifications are provided in Appx. A.6 and ablation studies in Sec. 4.

## 4. Experiments

### 4.1. Experimental Setup

**Physical constraints.** We adopt the unsafe-behavior taxonomy for bimanual manipulation from Deng et al. (2025), which encompasses the majority of hazard patterns identified through extensive case studies. To evaluate PACT, we select three representative constraint functions (i.e., Poking, Alignment, and Rotation) along with corresponding tasks (e.g., pouring water). The constraints require privileged environment states, which can be obtained directly through simulation APIs. For real-world experiments, we estimate these states using off-the-shelf 3D perception pipelines (Huang et al., 2025). See Appx. B for details.

### 4.2. Simulation Evaluation

**Environments and tasks.** We evaluate these methods on bimanual manipulation tasks from RoboTwin (Mu et al., 2024; Chen et al., 2025) with a *randomized* setup, requiring bimanual coordination and robust handling of safety-critical interactions. See Appx. C.1 & C.2 for details.

**Base policies and implementations.** To demonstrate generality, we apply PACT to multiple pretrained diffusion-based policies with a wide spectrum of model architectures, base

capacities, input modalities, diffusion parameterizations, and sampling strategies: Diffusion Policy (DP) (Chi et al., 2025), 3D Diffusion Policy (DP3) (Ze et al., 2024), which takes point clouds as model input, and generalist policies like RDT-1B (Liu et al., 2025b) and  $\pi_{0.5}$  (Black et al., 2025). During post-training, each policy collects data by rolling out on 1,000 training scenes. We perform full-parameter fine-tuning for DP and DP3, and apply Low-Rank Adaptation (LoRA) (Hu et al., 2022) to align RDT-1B and  $\pi_{0.5}$ . Additional details are provided in Appx. C.5.

**Baselines.** We organize baselines along two axes: learning paradigm (IL vs. RL) and data regime (off-policy vs. on-policy). Off-policy IL baselines include behavior cloning on intervened rollouts as in Probe, Learn, Distill (PLD) (Xiao et al., 2025) and cloning on rollouts collected with the implicit safe teacher (guided rollouts), while off-policy RL is represented by iDQL (Hansen-Estruch et al., 2023). To enable comparison between distillation and other objectives, we also evaluate offline distillation (Meng et al., 2023) with different types of rollouts. On-policy baselines use self-collected data and include IL methods (Rejection Fine-tuning (RFT) and online variant of PLD (PLD<sub>online</sub>)) and RL methods including AWR (Peng et al., 2019), QSM (Psenka et al., 2024), DIPO (Yang et al., 2023), and PPO-style diffusion optimization (Schulman et al., 2017; Liu et al., 2025a). All methods are initialized from the same pretrained DP and are matched in environment interaction and update budgets. Details are provided in Appx. C.6 and Appx. C.7.

**Metrics.** We report *Success Rate*, measuring task completion, and *Safe Rate*, measuring the fraction of rollouts that satisfy *all* physical safety demands. The details on the metrics computation are elaborated in Appx. C.4.

**Effectiveness of PACT.** The aggregate performance across tasks is summarized in Table 1. Overall, post-training with PACT consistently improves both task completion and manipulation safety across all base policies. For instance, DP shows substantial gains in both success rate by 28% and safety rate by 24%. It indicates that the aligned behaviors become simultaneously more effective and more

Table 2. Performance comparison with on-policy baselines. Both the Success Rate (Succ.) and Safe Rate (Safe) are reported.

Method	Pick Dual Bottles		Handover Apple		Pour Water		Stack Blocks	
	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe
<i>Imitation Learning</i>								
Base	67%	36%	63%	64%	38%	20%	63%	24%
PLD <sub>online</sub>	93%	66%	70%	70%	55%	29%	71%	32%
<i>Reinforcement Learning</i>								
PPO	71%	37%	67%	67%	60%	31%	73%	33%
DIPO	77%	43%	65%	68%	67%	29%	66%	36%
<i>Distillation</i>								
Ours	<b>96%</b>	<b>89%</b>	<b>81%</b>	<b>82%</b>	<b>86%</b>	<b>64%</b>	<b>80%</b>	<b>48%</b>

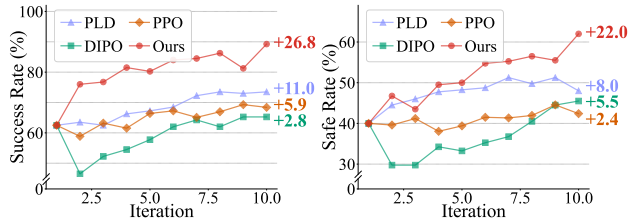


Figure 4. Training efficiency comparison with on-policy baselines. Success Rate (left) and Safe Rate (right), are averaged over four tasks across training iterations. Our method demonstrates the most training efficiency and stability.

physically compliant, which facilitates the reconciliation between performance and safety (Zhang et al., 2025a) through curriculum-based constraint distillation, which injects dense constraint gradients into the diffusion policy while progressively tightening safety enforcement to preserve behavioral fidelity. Furthermore, improvements are also more pronounced on difficult, precision-critical tasks (e.g., Pour Water and Stack Blocks), where minor misplacement can readily trigger constraint violations that cause failure. Qualitative results for simulated tasks are provided in Appx. E.1.

**Comparison to on-policy baselines.** We further compare PACT with representative on-policy baselines in Table 2. Across all four tasks, PACT attains the highest Success Rate and Safe Rate, indicating improved constraint satisfaction without sacrificing task performance. Beyond effectiveness, PACT is markedly more stable: several on-policy alternatives (e.g., AWR, RFT, and QSM) inevitably collapse to near-zero success in our setting (See Appx. E.3 for details), and RL-based baselines often require additional stabilization mechanisms (e.g., soft target updates in DIPO or KL regularization in PPO). In contrast, PACT leverages direct supervision from gradients of differentiable safety costs, avoiding the high-variance credit assignment inherent to policy-gradient updates whose quality is dominated by the sampled batch (Nota & Thomas, 2020). This stability is further reflected in Fig. 4, where PACT improves smoothly without the early drops observed in DIPO or the oscillations of PPO. Finally, PACT is the most efficient, converging more than  $5\times$  faster because constraint distillation provides dense, per-timestep supervision throughout the diffusion process. RL-based methods instead rely on coarse trajectory-level signals and costly, noisy likelihood estimation, and therefore

Table 3. Performance comparison with off-policy baselines. Success Rate (Succ.) and Safe Rate (Safe) are reported for each task. “RO” denotes abbreviation for Rollout.

Method	Pick Dual Bottles		Handover Apple		Pour Water		Stack Blocks	
	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe
<i>Imitation Learning</i>								
Base	67%	36%	63%	64%	38%	20%	63%	24%
Guided RO	82%	74%	74%	79%	70%	41%	66%	<b>59%</b>
PLD	70%	43%	70%	72%	54%	47%	71%	40%
<i>Reinforcement Learning</i>								
iDQL	74%	40%	63%	70%	62%	28%	65%	31%
<i>Distillation</i>								
Expert RO	70%	41%	71%	68%	67%	24%	65%	27%
Self RO	76%	44%	70%	67%	71%	28%	72%	34%
Guided RO	70%	45%	72%	78%	72%	31%	77%	28%
Ours	<b>96%</b>	<b>89%</b>	<b>81%</b>	<b>82%</b>	<b>86%</b>	<b>64%</b>	<b>80%</b>	48%

require substantially more interaction to reach comparable gains. Additional results are provided in Appx. E.

**Comparison to off-policy baselines.** We compare PACT against off-policy alternatives spanning imitation-based approaches (Guided RO and PLD) and an RL-based method (iDQL). As shown in Table 3, PACT attains higher success and safety rates in most settings, with particularly large gains on *Pick Dual Bottles* and *Pour Water*. We also implement an off-policy distillation variant under different data sources. Distillation is especially effective on more complex tasks (e.g., *Pour Water* and *Stack Blocks*), where collecting sufficient rollouts for imitation is difficult. However, distribution tilting via distillation can leverage limited data more effectively. Moreover, off-policy variants of PACT underperform the on-policy iterative update (ours) under matched gradient budgets, underscoring the importance of on-policy with the reverse-KL objective in Eq. (5), which mitigates forgetting in utility-free circumstances (Shenfeld et al., 2025) and avoids the constraint leakage induced by the mode-covering behavior of forward-KL.

**Ablation Studies. 1) Curriculum distillation.** We first ablate the curriculum used to progressively tighten safety constraints during training. Specifically, we compare direct distillation (Eq. (9)) to curriculum distillation with a linear schedule  $\eta_i = i/N$ . We find that direct distillation induces catastrophic forgetting with significant performance drops (Fig. 12) in contrast with the smoother improvement dynamics yielded by curriculum-based enforcement (Fig. 4). **2) Impact of Lagrange multipliers  $\lambda$ .** Next, a grid search in Fig. 6 reveals a trade-off between constraint enforcement and policy fidelity: overly large  $\lambda$  leads to irreversible OOD failure similar to direct distillation, whereas overly small  $\lambda$  yields slow and limited improvements. The best performance consistently occurs near a turning point, and PACT remains robust within a practical multiplier range. **3) Efficient distillation with few diffusion steps  $t_c$ .** As shown in Fig. 7, injecting constraints only for a small number of late diffusion steps achieves strong alignment while maintaining

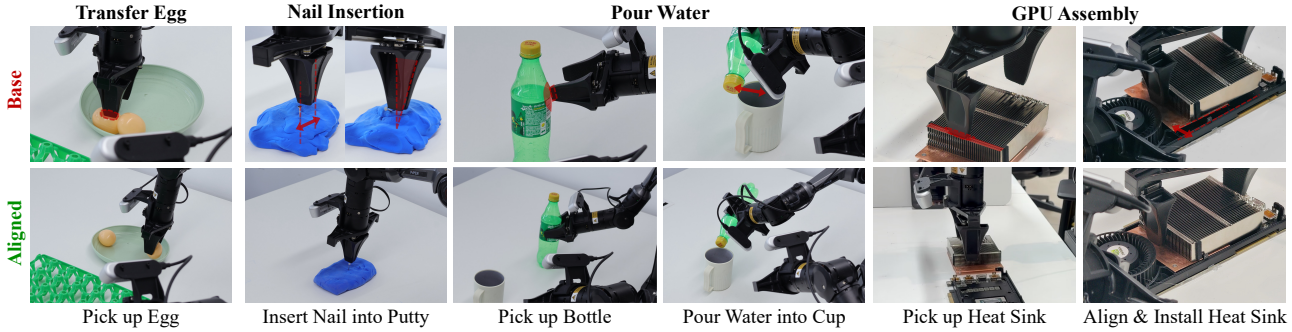


Figure 5. **Qualitative results of real world evaluation.** base policy (top) v.s. policy aligned by PACT (bottom) across four manipulation tasks. PACT reduces unsafe contacts and improves task completion by correcting key failure modes: avoiding poking to securely grasp the egg (*Transfer Egg*); aligning the gripper with the nail head, preventing lateral or tilted insertion (*Nail Insertion*); eliminating bottle poking and cup-rim misalignment (*Pour Water*); avoiding poking heat-sink and corrects installation misalignment (*GPU Assembly*).

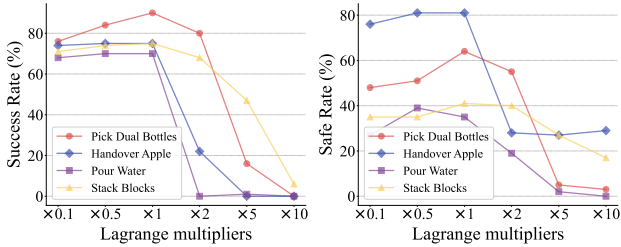


Figure 6. **Ablation study on Lagrange multipliers.** We report the performance of PACT after 5 iterations to reflect early-stage convergence behavior.

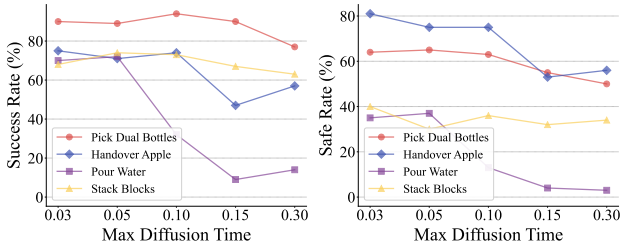


Figure 7. **Ablation study on max diffusion time for cost guidance  $t_c$ .** We report the performance of PACT after 5 iterations.

the approximation fidelity in Sec. 3.3, whereas increasing  $t_c$  often degrades performance and incurs additional computation, supporting our design of few-steps constraint distillation. **4) Sensitivity to rollouts and UTD ratios.** Finally, Table 4 shows that performance is relatively stable across a wide range of rollout counts and UTD ratios (number of epochs  $E$ ), indicating that PACT does not rely on excessive sampling or aggressive inner-loop optimization, but benefits from the dense and stable supervision by constraint distillation. More details are elaborated in Appx. E.5.

### 4.3. Real-World Evaluation

We evaluate PACT on safety-critical real-world manipulation tasks, including *Pour Water*, *Nail Insertion*, *Transfer Egg*, and *GPU Assembly*, and index corresponding constraints as in simulation (refer to Fig. 10 & Appx. D.2 for more details). *GPU Assembly* is particularly challenging, requiring millimeter-level precision to align heat sink position-

Table 4. **Effect of rollout count and update-to-data (UTD) ratio.** We report Success Rate and Safe Rate on each task after 5 iterations. † marks the default hyperparameters for DP.

Rollouts/ UTD Ratio	Pick Dual Bottles		Handover Apple		Pour Water		Stack Blocks	
	Succ.	Safe	Succ.	Safe	Succ.	Safe	Succ.	Safe
144 / 100	80%	48%	47%	51%	62%	29%	67%	38%
288 / 50	83%	51%	60%	64%	67%	15%	66%	30%
288 / 100†	90%	64%	75%	81%	70%	35%	68%	40%
288 / 200	91%	67%	45%	55%	71%	45%	69%	44%

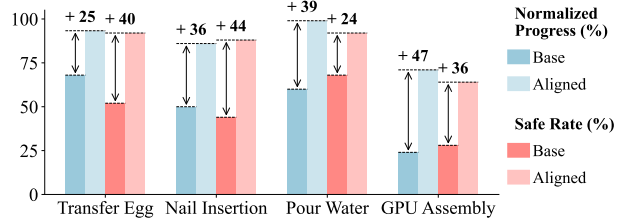


Figure 8. **Quantitative results of real-world evaluation.** We report normalized task progress and safe rate for the base policy vs. our aligned policy. PACT improves both metrics across all tasks, with the largest metric gain observed on GPU Assembly.

ing pins and mounting holes (Fig. 1). We use a scaled-down, instruction-free RDT (Liu et al., 2025b) variant as the base policy, trained from scratch on tele-operated demonstrations for each task (Appx. D). PACT substantially improves the physical Safe Rate by 36.0% and the average task progress by 36.8% across tasks as shown in Fig. 8. Furthermore, qualitatively (Fig. 5), PACT induces fine-grained, constraint-driven corrections that avoid hazardous interactions. These results validate post-training safety alignment for deploying diffusion-based policies in safety-critical robotic systems.

## 5. Conclusion

We introduced PACT, a post-training framework to align diffusion manipulation policies within the physical constraints. PACT distills cost gradients into the policy via a reverse-KL objective and applies a curriculum to mitigate catastrophic forgetting. Experiments show reduced safety violations and improved task success, supporting PACT as a practical approach for safety-critical embodied manipulation.

## Impact Statement

This paper studies post-training physical safety alignment for diffusion-based embodied manipulation policies, aiming to improve constraint compliance while preserving task capability. A potential positive impact is to enable safer deployment in safety-critical settings (e.g., manufacturing or assistive manipulation) by reducing hazardous behaviors without additional test-time guard modules or extensive expert supervision. However, stronger manipulation policies may increase dual-use risk, including misuse in harmful applications or deployment under misspecified or adversarially chosen constraints. We advocate responsible use by treating alignment as a risk-reduction mechanism rather than a formal safety guarantee.

## References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 20
- Ahmadian, A., Cremer, C., Gallé, M., Fadaee, M., Kretutzer, J., Pietquin, O., Üstün, A., and Hooker, S. Back to basics: Revisiting reinforce-style optimization for learning from human feedback in llms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12248–12267, 2024. 26
- Altman, E. *Constrained Markov decision processes*. Routledge, 2021. 2, 3
- Amin, A., Aniceto, R., Balakrishna, A., Black, K., Conley, K., Connors, G., Darpinian, J., Dhabalia, K., DiCarlo, J., et al.  $\pi_{0,6}^*$ : a vla that learns from experience. *arXiv preprint arXiv:2511.14759*, 2025. 1, 26, 27
- Bahety, A., Balaji, A., Abbatematteo, B., and Martín-Martín, R. Safemimic: Towards safe and autonomous human-to-robot imitation for mobile manipulation. *arXiv preprint arXiv:2506.15847*, 2025. 1, 3
- Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J., et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 30
- Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., Das-Sarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022. 2
- Ball, P. J., Smith, L., Kostrikov, I., and Levine, S. Efficient online reinforcement learning with offline data. In *International Conference on Machine Learning*, pp. 1577–1594. PMLR, 2023. 3
- Bao, F., Zhao, M., Hao, Z., Li, P., Li, C., and Zhu, J. Equivariant energy-guided sde for inverse molecular design. In *The Eleventh International Conference on Learning Representations*, 2022. 5, 19
- Bengio, S., Vinyals, O., Jaitly, N., and Shazeer, N. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28, 2015. 2
- Bi, H., Wu, L., Lin, T., Tan, H., Su, Z., Su, H., and Zhu, J. H-rdt: Human manipulation enhanced bimanual robotic manipulation. *arXiv preprint arXiv:2507.23523*, 2025. 30
- Billard, A. and Kragic, D. Trends and challenges in robot manipulation. *Science*, 364(6446):eaat8414, 2019. 2
- Black, K., Brown, N., Driess, D., Esmail, A., Equi, M., Finn, C., Fusai, N., Groom, L., Hausman, K., Ichter, B., et al.  $\pi_0$ : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024a. 16, 28
- Black, K., Janner, M., Du, Y., Kostrikov, I., and Levine, S. Training diffusion models with reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024b. 4
- Black, K., Brown, N., Darpinian, J., Dhabalia, K., Driess, D., Esmail, A., Equi, M. R., Finn, C., Fusai, N., Galliker, M. Y., Ghosh, D., Groom, L., Hausman, K., Ichter, B., Jakubczak, S., Jones, T., Ke, L., LeBlanc, D., Levine, S., Li-Bell, A., Mothukuri, M., Nair, S., Pertsch, K., Ren, A. Z., Shi, L. X., Smith, L., Springenberg, J. T., Stachowicz, K., Tanner, J., Vuong, Q., Walke, H., Walling, A., Wang, H., Yu, L., and Zhilinsky, U.  $\pi_{0,5}$ : a vision-language-action model with open-world generalization. In *Proceedings of The 9th Conference on Robot Learning*, volume 305 of *Proceedings of Machine Learning Research*, pp. 17–40. PMLR, 27–30 Sep 2025. 6, 16, 25, 28
- Brunke, L., Greeff, M., Hall, A. W., Yuan, Z., Zhou, S., Panerati, J., and Schoellig, A. P. Safe learning in robotics: From learning-based control to safe reinforcement learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 5(1):411–444, 2022. 3
- Buhl, J. F., Grønhoj, R., Jørgensen, J. K., Mateus, G., Pinto, D., Sørensen, J. K., Bøgh, S., and Chrysostomou, D. A dual-arm collaborative robot system for the smart factories of the future. *Procedia manufacturing*, 38:333–340, 2019. 1

- 495 Chen, H., Lu, C., Ying, C., Su, H., and Zhu, J. Offline rein-  
496 forcement learning via high-fidelity generative behavior  
497 modeling. In *The Eleventh International Conference on*  
498 *Learning Representations*, 2023. 2
- 500 Chen, H., Zheng, K., Su, H., and Zhu, J. Aligning diffusion  
501 behaviors with q-functions for efficient continuous con-  
502 trol. *Advances in Neural Information Processing Systems*,  
503 37:119949–119975, 2024. 2
- 504 Chen, T., Chen, Z., Chen, B., Cai, Z., Liu, Y., Li, Z., Liang,  
505 Q., Lin, X., Ge, Y., Gu, Z., et al. Robotwin 2.0: A scal-  
506 able data generator and benchmark with strong domain  
507 randomization for robust bimanual robotic manipulation.  
508 *arXiv preprint arXiv:2506.18088*, 2025. 6, 21, 23
- 510 Chi, C., Xu, Z., Feng, S., Cousineau, E., Du, Y., Burch-  
511 fiel, B., Tedrake, R., and Song, S. Diffusion policy:  
512 Visuomotor policy learning via action diffusion. *The*  
513 *International Journal of Robotics Research*, 44(10-11):  
514 1684–1704, 2025. 1, 2, 6, 23, 25
- 516 Chow, Y., Ghavamzadeh, M., Janson, L., and Pavone, M.  
517 Risk-constrained reinforcement learning with percentile  
518 risk criteria. *Journal of Machine Learning Research*, 18  
519 (167):1–51, 2018. 3
- 521 Chung, H., Kim, J., Mccann, M. T., Klasky, M. L., and  
522 Ye, J. C. Diffusion posterior sampling for general noisy  
523 inverse problems. In *The Tenth International Conference*  
524 *on Learning Representations*, 2023. 5, 6, 20
- 526 Ciftci, Y. U., Chiu, D., Feng, Z., Sukhatme, G. S., and  
527 Bansal, S. Safe-gil: Safety guided imitation learning  
528 for robotic systems. In *2025 IEEE International Confer-*  
529 *ence on Robotics and Automation (ICRA)*, pp. 3559–3566.  
530 IEEE, 2025. 1
- 532 Conn, A. R., Gould, N. I., and Toint, P. L. *Trust region*  
533 *methods*. Society for Industrial and Applied Mathematics  
534 (SIAM), 2000. 4
- 536 Dai, J., Chen, T., Wang, X., Yang, Z., Chen, T., Ji, J.,  
537 and Yang, Y. Safesora: Towards safety alignment of  
538 text2video generation via a human preference dataset.  
539 *Advances in Neural Information Processing Systems*, 37:  
540 17161–17214, 2024. 2
- 542 Dayan, P. and Hinton, G. E. Using expectation-  
543 maximization for reinforcement learning. *Neural Compu-*  
544 *tation*, 9(2):271–278, 1997. 3
- 546 Deng, H., Guo, W., Wang, Q., Wu, Z., and Wang, Z. Safebi-  
547 manual: Diffusion-based trajectory optimization for safe  
548 bimanual manipulation. In *Conference on Robot Learn-*  
549 *ing*, pp. 3218–3238. PMLR, 2025. 1, 3, 6, 20, 21
- Dulac-Arnold, G., Levine, N., Mankowitz, D. J., Li, J.,  
Paduraru, C., Gowal, S., and Hester, T. Challenges of  
real-world reinforcement learning: definitions, bench-  
marks and analysis. *Machine Learning*, 110(9):2419–  
2468, 2021. 3
- Fu, Z., Zhao, T. Z., and Finn, C. Mobile aloha: Learning bi-  
manual mobile manipulation using low-cost whole-body  
teleoperation. In *Conference on Robot Learning*, pp.  
4066–4083. PMLR, 2025. 27
- Garcia, J. and Fernández, F. A comprehensive survey on safe  
reinforcement learning. *Journal of Machine Learning*  
*Research*, 16(1):1437–1480, 2015. 1, 2, 3
- Gilks, W. R. and Wild, P. Adaptive rejection sampling for  
gibbs sampling. *Journal of the Royal Statistical Society:*  
*Series C (Applied Statistics)*, 41(2):337–348, 1992. 31
- Gokmen, C., Ho, D., and Khansari, M. Asking for help:  
Failure prediction in behavioral cloning through value ap-  
proximation. In *2023 IEEE International Conference on*  
*Robotics and Automation (ICRA)*, pp. 5821–5828. IEEE,  
2023. 3
- Gu, S., Yang, L., Du, Y., Chen, G., Walter, F., Wang, J.,  
and Knoll, A. A review of safe reinforcement learning:  
Methods, theories, and applications. *IEEE Transactions*  
*on Pattern Analysis and Machine Intelligence*, 46(12):  
11216–11235, 2024. 3
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft  
actor-critic: Off-policy maximum entropy deep reinforce-  
ment learning with a stochastic actor. In *International*  
*conference on machine learning*, pp. 1861–1870. Pmlr,  
2018. 31
- Haddadin, S. Physical safety in robotics. In *Formal Mod-*  
*eling and Verification of Cyber-Physical Systems: 1st*  
*International Summer School on Methods and Tools for*  
*the Design of Digital Systems, Bremen, Germany, Septem-*  
*ber 2015*, pp. 249–271. Springer, 2015. 3
- Hansen-Estruch, P., Kostrikov, I., Janner, M., Kuba, J. G.,  
and Levine, S. Idql: Implicit q-learning as an actor-  
critic method with diffusion policies. *arXiv preprint*  
*arXiv:2304.10573*, 2023. 6, 26
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-  
bilistic models. *Advances in neural information process-*  
*ing systems*, 33:6840–6851, 2020. 2, 23
- Hsu, K.-C., Hu, H., and Fisac, J. F. The safety filter: A uni-  
fied view of safety-critical control in autonomous systems.  
*Annual Review of Control, Robotics, and Autonomous*  
*Systems*, 7, 2023. 1

- 550 Hu, E. J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang,  
551 L., Chen, W., et al. Lora: Low-rank adaptation of large  
552 language models. In *The Tenth International Conference*  
553 *on Learning Representations*, 2022. 6, 23
- 554 Hu, Z. J., Wang, Z., Huang, Y., Sena, A., y Baena, F. R., and  
555 Burdet, E. Towards human-robot collaborative surgery:  
556 Trajectory and strategy learning in bimanual peg transfer.  
557 *IEEE Robotics and Automation Letters*, 8(8):4553–4560,  
558 2023. 1
- 560 Huang, W., Wang, C., Zhang, R., Li, Y., Wu, J., and Fei-Fei,  
561 L. Voxposer: Composable 3d value maps for robotic  
562 manipulation with language models. *Proceedings of Ma-*  
563 *chine Learning Research*, 229, 2023. 20
- 565 Huang, W., Wang, C., Li, Y., Zhang, R., and Fei-Fei, L.  
566 Rekep: Spatio-temporal reasoning of relational keypoint  
567 constraints for robotic manipulation. In *Conference on*  
568 *Robot Learning*, pp. 4573–4602. PMLR, 2025. 6, 20
- 569 Janner, M., Du, Y., Tenenbaum, J., and Levine, S. Planning  
570 with diffusion for flexible behavior synthesis. In *Internat-*  
571 *ional Conference on Machine Learning*, pp. 9902–9915.  
572 PMLR, 2022. 1, 2
- 574 Kalakrishnan, M., Chitta, S., Theodorou, E., Pastor, P., and  
575 Schaal, S. Stomp: Stochastic trajectory optimization for  
576 motion planning. In *2011 IEEE international conference*  
577 *on robotics and automation*, pp. 4569–4574. IEEE, 2011.  
578 20
- 580 Kingma, D. and Gao, R. Understanding diffusion objectives  
581 as the elbo with simple data augmentation. *Advances*  
582 *in Neural Information Processing Systems*, 36:65484–  
583 65516, 2023. 3
- 584 Lee, J., Paduraru, C., Mankowitz, D. J., Heess, N., Pre-  
585 cup, D., Kim, K.-E., and Guez, A. Coptidice: Offline  
586 constrained reinforcement learning via stationary distri-  
587 bution correction estimation. In *The Tenth International*  
588 *Conference on Learning Representations*, 2022. 3
- 590 Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and  
591 Le, M. Flow matching for generative modeling. In *The*  
592 *Eleventh International Conference on Learning Re-*  
593 *presentations*, 2023. 3, 15
- 594 Liu, H., Dass, S., Martín-Martín, R., and Zhu, Y. Model-  
595 based runtime monitoring with interactive imitation learn-  
596 ing. In *2024 IEEE International Conference on Robotics*  
597 *and Automation (ICRA)*, pp. 4154–4161. IEEE, 2024. 1,  
598 2, 3, 16
- 600 Liu, J., Liu, G., Liang, J., Li, Y., Liu, J., Wang, X., Wan,  
601 P., Zhang, D., and Ouyang, W. Flow-grpo: Training  
602 flow matching models via online rl. *arXiv preprint*  
603 *arXiv:2505.05470*, 2025a. 4, 6, 27
- Liu, S., Wu, L., Li, B., Tan, H., Chen, H., Wang, Z., Xu,  
K., Su, H., and Zhu, J. Rdt-1b: a diffusion foundation  
model for bimanual manipulation. In *The Thirteenth*  
*International Conference on Learning Representations*,  
2025b. 1, 6, 8, 23, 25
- Liu, X., Gong, C., et al. Flow straight and fast: Learning  
to generate and transfer data with rectified flow. In *The*  
*Eleventh International Conference on Learning Re-*  
*presentations*, 2023. 3, 15, 23, 30
- Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., and Zhu, J.  
Dpm-solver: A fast ode solver for diffusion probabilistic  
model sampling in around 10 steps. *Advances in neural*  
*information processing systems*, 35:5775–5787, 2022. 23
- Lu, C., Chen, H., Chen, J., Su, H., Li, C., and Zhu, J. Con-  
trastive energy prediction for exact energy-guided diffu-  
sion sampling in offline reinforcement learning. In *Inter-*  
*national Conference on Machine Learning*, pp. 22825–  
22855. PMLR, 2023. 4, 5, 14, 15
- Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., and Zhu, J. Dpm-  
solver++: Fast solver for guided sampling of diffusion  
probabilistic models. *Machine Intelligence Research*, pp.  
1–22, 2025. 23
- McAllister, D., Ge, S., Yi, B., Kim, C. M., Weber, E., Choi,  
H., Feng, H., and Kanazawa, A. Flow matching policy  
gradients. *arXiv preprint arXiv:2507.21053*, 2025. 27
- McCloskey, M. and Cohen, N. J. Catastrophic interfer-  
ence in connectionist networks: The sequential learning  
problem. In *Psychology of learning and motivation*, vol-  
ume 24, pp. 109–165. Elsevier, 1989. 2
- Meng, C., Rombach, R., Gao, R., Kingma, D., Ermon, S.,  
Ho, J., and Salimans, T. On distillation of guided diffu-  
sion models. In *Proceedings of the IEEE/CVF conference*  
*on computer vision and pattern recognition*, pp. 14297–  
14306, 2023. 6, 24
- Mu, Y., Chen, T., Peng, S., Chen, Z., Gao, Z., Zou, Y.,  
Lin, L., Xie, Z., and Luo, P. Robotwin: Dual-arm robot  
benchmark with generative digital twins (early version).  
In *European Conference on Computer Vision*, pp. 264–  
273. Springer, 2024. 6, 21, 23
- Nota, C. and Thomas, P. S. Is the policy gradient a gradient?  
In *Proceedings of the 19th International Conference on*  
*Autonomous Agents and MultiAgent Systems*, pp. 939–  
947, 2020. 7
- Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec,  
M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-  
Nouby, A., et al. Dinov2: Learning robust visual features  
without supervision. *Transactions on Machine Learning*  
*Research Journal*, 2024. 20

- 605 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J.,  
 606 Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga,  
 607 L., et al. Pytorch: An imperative style, high-performance  
 608 deep learning library. *Advances in neural information  
 609 processing systems*, 32, 2019. 23
- 610 Pearce, T., Rashid, T., Kanervisto, A., Bignell, D., Sun, M.,  
 611 Georgescu, R., Macua, S. V., Tan, S. Z., Momennejad, I.,  
 612 Hofmann, K., et al. Imitating human behaviour with dif-  
 613 fusion models. In *The Eleventh International Conference  
 614 on Learning Representations*, 2023. 1
- 615 Peng, X. B., Kumar, A., Zhang, G., and Levine, S.  
 616 Advantage-weighted regression: Simple and scalable  
 617 off-policy reinforcement learning. *arXiv preprint  
 618 arXiv:1910.00177*, 2019. 3, 6, 26, 31
- 619 Psenka, M., Escontrela, A., Abbeel, P., and Ma, Y. Learning  
 620 a diffusion model policy from rewards via q-score match-  
 621 ing. In *Proceedings of the 41st International Conference  
 622 on Machine Learning*, pp. 41163–41182, 2024. 6, 26, 31
- 623 Ravi, N., Gabeur, V., Hu, Y.-T., Hu, R., Ryali, C., Ma, T.,  
 624 Khedr, H., Rädle, R., Rolland, C., Gustafson, L., et al.  
 625 Sam 2: Segment anything in images and videos. In  
 626 *The Thirteenth International Conference on Learning  
 627 Representations*, 2025. 20
- 628 Reichlin, A., Marchetti, G. L., Yin, H., Ghadirzadeh, A., and  
 629 Kragic, D. Back to the manifold: Recovering from out-  
 630 of-distribution states. In *2022 IEEE/RSJ International  
 631 Conference on Intelligent Robots and Systems (IROS)*, pp.  
 632 8660–8666. IEEE, 2022. 1, 3
- 633 Ren, A. Z., Lidard, J., Ankile, L. L., Simeonov, A., Agrawal,  
 634 P., Majumdar, A., Burchfiel, B., Dai, H., and Simchowitz,  
 635 M. Diffusion policy optimization. *arXiv preprint  
 636 arXiv:2409.00588*, 2024. 27
- 637 Ross, S., Gordon, G., and Bagnell, D. A reduction of imita-  
 638 tion learning and structured prediction to no-regret online  
 639 learning. In *Proceedings of the fourteenth international  
 640 conference on artificial intelligence and statistics*, pp.  
 641 627–635. JMLR Workshop and Conference Proceedings,  
 642 2011. 1, 3, 26
- 643 Rueß, H. Systems challenges for trustworthy embodied  
 644 systems. *arXiv preprint arXiv:2201.03413*, 2022. 1
- 645 Schulman, J., Duan, Y., Ho, J., Lee, A., Awwal, I., Bradlow,  
 646 H., Pan, J., Patil, S., Goldberg, K., and Abbeel, P. Motion  
 647 planning with sequential convex optimization and convex  
 648 collision checking. *The International Journal of Robotics  
 649 Research*, 33(9):1251–1270, 2014. 20
- 650 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and  
 651 Klimov, O. Proximal policy optimization algorithms.  
 652 *arXiv preprint arXiv:1707.06347*, 2017. 3, 6, 27, 31
- 653 Shenfeld, I., Pari, J., and Agrawal, P. RL’s razor: Why  
 654 online reinforcement learning forgets less. *arXiv preprint  
 655 arXiv:2509.04259*, 2025. 7
- 656 Siméoni, O., Vo, H. V., Seitzer, M., Baldassarre, F., Oquab,  
 657 M., Jose, C., Khalidov, V., Szafraniec, M., Yi, S.,  
 658 Ramamonjisoa, M., et al. Dinov3. *arXiv preprint  
 659 arXiv:2508.10104*, 2025. 30
- 660 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and  
 661 Ganguli, S. Deep unsupervised learning using nonequi-  
 662 librium thermodynamics. In *International conference on  
 663 machine learning*, pp. 2256–2265. pmlr, 2015. 2
- 664 Song, J., Meng, C., and Ermon, S. Denoising diffusion  
 665 implicit models. In *The Ninth International Conference  
 666 on Learning Representations*, 2021a. 3, 23
- 667 Song, Y., Durkan, C., Murray, I., and Ermon, S. Maxi-  
 668 mum likelihood training of score-based diffusion models.  
 669 *Advances in neural information processing systems*, 34:  
 670 1415–1428, 2021b. 3
- 671 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Er-  
 672 mon, S., and Poole, B. Score-based generative modeling  
 673 through stochastic differential equations. In *The Ninth  
 674 International Conference on Learning Representations*,  
 675 2021c. 2
- 676 Tessler, C., Mankowitz, D. J., and Mannor, S. Re-  
 677 ward constrained policy optimization. *arXiv preprint  
 678 arXiv:1805.11074*, 2018. 3
- 679 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,  
 680 L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. At-  
 681 tention is all you need. *Advances in neural information  
 682 processing systems*, 30, 2017. 23
- 683 Wabersich, K. P., Taylor, A. J., Choi, J. J., Sreenath, K.,  
 684 Tomlin, C. J., Ames, A. D., and Zeilinger, M. N. Data-  
 685 driven safety filters: Hamilton-jacobi reachability, control  
 686 barrier functions, and predictive methods for uncertain  
 687 systems. *IEEE Control Systems Magazine*, 43(5):137–  
 688 177, 2023. 3
- 689 Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi,  
 690 E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting  
 691 elicits reasoning in large language models. *Advances in  
 692 neural information processing systems*, 35:24824–24837,  
 693 2022. 20
- 694 Williams, R. J. Simple statistical gradient-following algo-  
 695 rithms for connectionist reinforcement learning. *Machine  
 696 learning*, 8(3):229–256, 1992. 26

- 660 Wong, J., Tung, A., Kurenkov, A., Mandlkar, A., Fei-Fei,  
661 L., Savarese, S., and Martín-Martín, R. Error-aware imita-  
662 tion learning from teleoperation data for mobile manipu-  
663 lation. In *Conference on Robot Learning*, pp. 1367–1378.  
664 PMLR, 2022. 1, 3
- 665 Xiao, W., Lin, H., Peng, A., Xue, H., He, T., Xie, Y., Hu, F.,  
666 Wu, J., Luo, Z., Fan, L., Shi, G., and Zhu, Y. Probe, learn,  
667 distill: Self-improving vision-language-action models  
668 with data generation via residual rl. *arXiv preprint*, 2025.  
669 6, 21, 24, 26
- 670 Xie, Z., Zhang, Q., Yang, F., Hutter, M., and Xu, R. Simple  
671 policy optimization. In *Proceedings of the 42nd Inter-  
672 national Conference on Machine Learning*, pp. 68813–  
673 68824, 2024. 27
- 674 Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang,  
675 J., and Dong, Y. Imagereward: Learning and evaluat-  
676 ing human preferences for text-to-image generation. *Ad-  
677 vances in Neural Information Processing Systems*, 36:  
678 15903–15935, 2023. 5
- 679 Yang, L., Huang, Z., Lei, F., Zhong, Y., Yang, Y., Fang, C.,  
680 Wen, S., Zhou, B., and Lin, Z. Policy representation via  
681 diffusion probability model for reinforcement learning.  
682 *arXiv preprint arXiv:2305.13122*, 2023. 6, 27
- 683 Yang, Y., Chen, L., Zaidi, Z., van Waveren, S., Krishna,  
684 A., and Gombolay, M. Enhancing safety in learning  
685 from demonstration algorithms via control barrier func-  
686 tion shielding. In *Proceedings of the 2024 ACM/IEEE  
687 International Conference on Human-Robot Interaction*,  
688 pp. 820–829, 2024. 3
- 689 Ye, H., Lin, H., Han, J., Xu, M., Liu, S., Liang, Y., Ma,  
690 J., Zou, J. Y., and Ermon, S. Tfg: Unified training-free  
691 guidance for diffusion models. *Advances in Neural In-  
692 formation Processing Systems*, 37:22370–22417, 2024.  
693 24
- 694 Ying, C., Zhou, X., Su, H., Yan, D., Chen, N., and Zhu, J.  
695 Towards safe reinforcement learning via constraining con-  
696 ditional value-at-risk. In *Proceedings of the Thirty-First  
697 International Joint Conference on Artificial Intelligence*,  
698 pp. 3673–3680, 2022. 3
- 699 Ying, C., Chen, H., Zhou, X., Hao, Z., Su, H., and Zhu, J.  
700 Exploratory diffusion model for unsupervised reinforce-  
701 ment learning. *arXiv preprint arXiv:2502.07279*, 2025.  
702 24
- 703 Yu, J., Wang, Y., Zhao, C., Ghanem, B., and Zhang, J. Free-  
704 dom: Training-free energy-guided conditional diffusion  
705 model. In *Proceedings of the IEEE/CVF International  
706 Conference on Computer Vision*, pp. 23174–23184, 2023.  
707 6
- 708 Ze, Y., Zhang, G., Zhang, K., Hu, C., Wang, M., and Xu,  
709 H. 3d diffusion policy: Generalizable visuomotor policy  
710 learning via simple 3d representations. In *ICRA 2024  
711 Workshop on 3D Visual Representations for Robot Ma-  
712 nipulation*, 2024. 6, 23, 25
- 713 Zhai, A., Liu, B., Fang, B., Cai, C., Ma, E., Yin, E., Wang,  
714 H., Zhou, H., Wang, J., Shi, L., et al. Igniting vlms toward  
the embodied space. *arXiv preprint arXiv:2509.11766*,  
2025. 16
- Zhang, B., Zhang, Y., Ji, J., Lei, Y., Dai, J., Chen, Y., and  
Yang, Y. Safevla: Towards safety alignment of vision-  
language-action model via constrained learning. *arXiv  
preprint arXiv:2503.03480*, 2025a. 3, 7
- Zhang, J., Huang, W., Peng, B., Wu, M., Hu, F., Chen, Z.,  
Zhao, B., and Dong, H. Omni6dpose: A benchmark  
and model for universal 6d object pose estimation and  
tracking. In *European Conference on Computer Vision*.  
Springer, 2024. 20
- Zhang, Y., Zhang, S., Huang, Y., Xia, Z., Fang, Z., Yang,  
X., Duan, R., Yan, D., Dong, Y., and Zhu, J. Stair: Im-  
proving safety alignment with introspective reasoning.  
In *International Conference on Machine Learning*, pp.  
76754–76777. PMLR, 2025b. 2
- Zhao, M., Bao, F., Li, C., and Zhu, J. Egsde: Unpaired  
image-to-image translation via energy-guided stochastic  
differential equations. *Advances in Neural Information  
Processing Systems*, 35:3609–3623, 2022. 19
- Zhao, T., Kumar, V., Levine, S., and Finn, C. Learning fine-  
grained bimanual manipulation with low-cost hardware.  
*Robotics: Science and Systems XIX*, 2023. 21
- Zheng, K., Lu, C., Chen, J., and Zhu, J. Improved techniques  
for maximum likelihood estimation for diffusion odes.  
In *International Conference on Machine Learning*, pp.  
42363–42389. PMLR, 2023. 15
- Zheng, Y., Li, J., Yu, D., Yang, Y., Li, S. E., Zhan, X., and  
Liu, J. Safe offline reinforcement learning with feasibility-  
guided diffusion model. In *The Twelfth International  
Conference on Learning Representations*, 2024. 26

## A. Proofs and Additional Theory

### A.1. Intermediate Cost

Following Lu et al. (2023), we formally analyze the structure of intermediate cost  $c_{k,t}(\mathbf{a}_t, \mathbf{s})$ . By rewriting  $\pi_0^* \triangleq \pi^*$  and  $\mu_{\phi,0} \triangleq \mu_\phi$ , the reformatted constrained optimal policy in Eq. (7) is:

$$\pi_0^*(\mathbf{a}_0 | \mathbf{s}) \propto \mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right),$$

Extending Lu et al. (2023), the exact intermediate cost is defined as follows.

**Theorem A.1** (Intermediate Cost Gradient). *For  $t \in (0, 1]$ , we have*

$$\mu_{\phi,t|0}(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \triangleq \pi_{t|0}^*(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) = \mathcal{N}(\mathbf{a}_t | \alpha_t \mathbf{a}_0, \sigma_t^2 \mathbf{I}).$$

Denote  $\pi_t^*(\mathbf{a}_t | \mathbf{s}) \triangleq \int \pi_{t|0}^*(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \pi_0^*(\mathbf{a}_0 | \mathbf{s}) d\mathbf{a}_0$  and  $\mu_{\phi,t}(\mathbf{a}_t | \mathbf{s}) \triangleq \int \mu_{\phi,t|0}(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s}) d\mathbf{a}_0$  as the marginal distributions at time  $t$ ; the Intermediate Cost is defined as

$$c_{k,t}(\mathbf{a}_t, \mathbf{s}) \triangleq \begin{cases} c_k(\mathbf{a}_0, \mathbf{s}), & t = 0, \\ -\frac{1}{\lambda_k} \log \mathbb{E}_{\pi_{0|t}^*(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})} \left[ \exp(-\lambda_k \cdot c_k(\mathbf{a}_0, \mathbf{s})) \right], & 0 < t \leq 1. \end{cases} \quad (16)$$

And their score functions satisfy

$$\nabla_{\mathbf{a}_t} \pi_t^*(\mathbf{a}_t | \mathbf{s}, t) = \underbrace{\nabla_{\mathbf{a}_t} \mu_\phi(\mathbf{a}_t | \mathbf{s}, t)}_{=\epsilon_\phi(\mathbf{a}_t, \mathbf{s}, t)/\sigma_t} - \underbrace{\sum_{k=1}^m \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s})}_{\text{intermediate cost gradient (intractable)}}. \quad (17)$$

*Proof.* We first note that the forward diffusion kernel is independent of the data distribution, hence for any  $t \in (0, 1]$ ,

$$\mu_{\phi,t|0}(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) = \pi_{t|0}^*(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) = \mathcal{N}(\mathbf{a}_t | \alpha_t \mathbf{a}_0, \sigma_t^2 \mathbf{I}).$$

Define the normalizing constant of the constrained optimum at  $t = 0$  as

$$Z(\mathbf{s}) \triangleq \int \mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right) d\mathbf{a}_0 = \mathbb{E}_{\mu_{\phi,0}(\cdot | \mathbf{s})} \left[ \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right) \right].$$

Then Eq. (7) can be written as

$$\pi_0^*(\mathbf{a}_0 | \mathbf{s}) = \frac{\mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right)}{Z(\mathbf{s})}.$$

For  $t \in (0, 1]$ , the marginal  $\pi_t^*(\mathbf{a}_t | \mathbf{s})$  is

$$\pi_t^*(\mathbf{a}_t | \mathbf{s}) = \int \pi_{t|0}^*(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \pi_0^*(\mathbf{a}_0 | \mathbf{s}) d\mathbf{a}_0 = \int \mu_{\phi,t|0}(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s}) \frac{\exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right)}{Z(\mathbf{s})} d\mathbf{a}_0.$$

Let

$$\mu_{\phi,t}(\mathbf{a}_t | \mathbf{s}) \triangleq \int \mu_{\phi,t|0}(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s}) d\mathbf{a}_0,$$

and

$$\mu_{\phi,0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s}) \triangleq \frac{\mu_{\phi,t|0}(\mathbf{a}_t | \mathbf{a}_0, \mathbf{s}) \mu_{\phi,0}(\mathbf{a}_0 | \mathbf{s})}{\mu_{\phi,t}(\mathbf{a}_t | \mathbf{s})}.$$

Substituting the posterior identity yields

$$\pi_t^*(\mathbf{a}_t | \mathbf{s}) = \mu_{\phi,t}(\mathbf{a}_t | \mathbf{s}) \frac{\mathbb{E}_{\mu_{\phi,0|t}(\cdot|\mathbf{a}_t,\mathbf{s})} \left[ \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right) \right]}{Z(\mathbf{s})}$$

To match the multiplicative form in Eq. (7), define the intermediate cost for each  $k$  by

$$c_{k,t}(\mathbf{a}_t, \mathbf{s}) \triangleq \begin{cases} c_k(\mathbf{s}, \mathbf{a}_0), & t = 0, \\ -\log \mathbb{E}_{\mu_{\phi,0|t}(\mathbf{a}_0|\mathbf{a}_t,\mathbf{s})} \left[ \exp\left(-\lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right) \right] / \lambda_k, & 0 < t \leq 1, \end{cases}$$

so that

$$\exp\left(-\sum_{k=1}^m \lambda_k c_{k,t}(\mathbf{a}_t, \mathbf{s})\right) = \prod_{k=1}^m \mathbb{E}_{\mu_{\phi,0|t}(\mathbf{a}_0|\mathbf{a}_t,\mathbf{s})} \left[ \exp\left(-\lambda_k c_k(\mathbf{s}, \mathbf{a}_0)\right) \right],$$

and therefore

$$\pi_t^*(\mathbf{a}_t | \mathbf{s}) \propto \mu_{\phi,t}(\mathbf{a}_t | \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_{k,t}(\mathbf{a}_t, \mathbf{s})\right).$$

Taking gradients with respect to  $\mathbf{a}_t$  gives the score decomposition

$$\nabla_{\mathbf{a}_t} \log \pi_t^*(\mathbf{a}_t | \mathbf{s}) = \nabla_{\mathbf{a}_t} \log \mu_{\phi,t}(\mathbf{a}_t | \mathbf{s}) - \sum_{k=1}^m \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s}),$$

which is exactly the claimed result.  $\square$

Following Lu et al. (2023), we view the *intermediate cost*  $c_{k,t}(\mathbf{a}_t, \mathbf{s})$  as the diffusion-time analogue of the original constraint  $c_k(\mathbf{a}_0, \mathbf{s})$  defined on clean actions. Although the constrained optimum is specified at  $t = 0$  in Eq. (7), the reverse diffusion process must operate on noisy actions  $\mathbf{a}_t$  for  $t > 0$ . Theorem A.1 shows that the marginal  $\pi_t^*(\mathbf{a}_t | \mathbf{s})$  retains the same multiplicative structure at each diffusion time, with the original constraint replaced by an intermediate cost  $c_{k,t}(\mathbf{a}_t, \mathbf{s})$ . In particular,  $c_{k,t}$  takes a log-expectation form under the posterior  $\pi_{t|0}^*(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})$ , and can be interpreted as a soft aggregation of the costs of plausible clean actions  $\mathbf{a}_0$  that could have produced the current noisy action  $\mathbf{a}_t$ .

This perspective also explains the score decomposition in Theorem A.1. The score of  $\pi_t^*$  decomposes into a base score term provided by the pretrained diffusion model and the *intermediate cost gradient*  $\nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s})$ . While this decomposition is exact, evaluating  $\nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s})$  is generally intractable due to the log-expectation over  $\mathbf{a}_0$  under the posterior induced by  $\pi_0^*$ . Consequently, principled constrained sampling from  $\pi^*$  requires approximating the intermediate cost guidance, which motivates practical approximations.

## A.2. Parameterization-Agnostic Attribute of Distillation Objective in Eq. (9)

Although our main distillation objective in Eq. (9) is written using the  $\epsilon$ -parameterization, we further show that the objective is *parameterization-agnostic*: the same teacher–student alignment can be expressed equivalently under alternative diffusion parameterizations (e.g.,  $v$ -parameterization (Zheng et al., 2023) and flow matching (Lipman et al., 2023; Liu et al., 2023)) by applying an invertible, time-dependent linear transform for each diffusion time  $t$ . Under the forward process notation

$$\mathbf{a}_t = \alpha_t \mathbf{a}_0 + \sigma_t \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

the  $v$ -parameterization is defined by the schedule derivatives as

$$\mathbf{v}_t \triangleq \dot{\alpha}_t \mathbf{a}_0 + \dot{\sigma}_t \boldsymbol{\epsilon},$$

where  $\dot{\alpha}_t \triangleq \frac{d\alpha_t}{dt}$  and  $\dot{\sigma}_t \triangleq \frac{d\sigma_t}{dt}$ . Stacking the two relations yields a linear system:

$$\begin{bmatrix} \mathbf{a}_t \\ \mathbf{v}_t \end{bmatrix} = \begin{bmatrix} \alpha_t & \sigma_t \\ \dot{\alpha}_t & \dot{\sigma}_t \end{bmatrix} \begin{bmatrix} \mathbf{a}_0 \\ \boldsymbol{\epsilon} \end{bmatrix}.$$

For any  $t \in (0, 1]$  such that the determinant

$$\Delta_t \triangleq \alpha_t \dot{\sigma}_t - \sigma_t \dot{\alpha}_t$$

is nonzero (which holds for standard noise schedules), the mapping is invertible. In particular,

$$\mathbf{a}_0 = \frac{\dot{\sigma}_t \mathbf{a}_t - \sigma_t \mathbf{v}_t}{\Delta_t}, \quad \boldsymbol{\epsilon} = \frac{-\dot{\alpha}_t \mathbf{a}_t + \alpha_t \mathbf{v}_t}{\Delta_t}.$$

Equivalently, given an  $\boldsymbol{\epsilon}$ -prediction, one can recover  $\mathbf{v}$  via

$$\mathbf{a}_0 = \frac{\mathbf{a}_t - \sigma_t \boldsymbol{\epsilon}}{\alpha_t}, \quad \mathbf{v} = \dot{\alpha}_t \mathbf{a}_0 + \dot{\sigma}_t \boldsymbol{\epsilon},$$

and similarly given  $\mathbf{v}$ -prediction one can recover  $\boldsymbol{\epsilon}$  via the inverse transform above. Therefore, any teacher signal expressed in  $\boldsymbol{\epsilon}$ -space can be converted to an equivalent teacher in  $\mathbf{v}$ -space and vice versa.

**Equivalent distillation objectives.** Let the implicit teacher be defined in  $\boldsymbol{\epsilon}$ -parameterization as  $\boldsymbol{\epsilon}^*(\mathbf{a}_t, \mathbf{s}, t)$ , and define the corresponding velocity teacher by

$$\mathbf{v}^*(\mathbf{a}_t, \mathbf{s}, t) \triangleq \dot{\alpha}_t \hat{\mathbf{a}}_0^*(\mathbf{a}_t, \mathbf{s}, t) + \dot{\sigma}_t \boldsymbol{\epsilon}^*(\mathbf{a}_t, \mathbf{s}, t), \quad \hat{\mathbf{a}}_0^*(\mathbf{a}_t, \mathbf{s}, t) \triangleq \frac{\mathbf{a}_t - \sigma_t \boldsymbol{\epsilon}^*(\mathbf{a}_t, \mathbf{s}, t)}{\alpha_t}.$$

Then the score-matching distillation objective in Eq. (9) can be equivalently written under  $\mathbf{v}$ -parameterization as

$$\min_{\theta} \mathbb{E}_{t, \boldsymbol{\epsilon}, (\mathbf{s}, \mathbf{a}) \sim d^{\pi_{\theta}}} \left[ \|\mathbf{v}_{\theta}(\mathbf{a}_t, \mathbf{s}, t) - \mathbf{v}^*(\mathbf{a}_t, \mathbf{s}, t)\|_2^2 \right],$$

where  $\mathbf{v}_{\theta}$  is the student policy expressed in  $\mathbf{v}$ -parameterization. Since the transform between  $(\mathbf{a}_0, \boldsymbol{\epsilon})$  and  $(\mathbf{a}_t, \mathbf{v}_t)$  is linear and invertible for each  $t$ , minimizing a squared error in one parameterization induces a corresponding squared error in the other (up to a deterministic, schedule-dependent reweighting). Consequently, our alignment objective is parameterization-agnostic: the same teacher distribution can be distilled into the student regardless of whether the diffusion model is implemented via  $\boldsymbol{\epsilon}$ -prediction,  $\mathbf{v}$ -prediction, or other equivalent parameterizations.

Moreover, the Rectified flow (Liu et al., 2024), which is widely used in the state-of-the-art Vision-Language-Action Model (VLA) with diffusion action expert (Black et al., 2024a; 2025; Zhai et al., 2025) can be viewed as a simplified special case of  $\mathbf{v}$ -parameterization with  $\alpha_t = 1 - t$ , and  $\sigma_t = t$ , resulting in the target velocity as  $\mathbf{v} = \boldsymbol{\epsilon} - \mathbf{a}_0$ .

### A.3. Proof of Theorem 3.1

**Lagrangian minimizer.** The following lemma states the closed-form pointwise minimizer of the KL-regularized Lagrangian objective from Eq. (6) and connects it directly to the Boltzmann-tilted policy in Eq. (7).

**Lemma A.2** (Pointwise minimizer of the KL-regularized Lagrangian). *Fix a state  $\mathbf{s}$  and multipliers  $\lambda_k \geq 0$ . Consider*

$$\mathcal{F}(\pi(\cdot | \mathbf{s})) \triangleq D_{\text{KL}}(\pi(\cdot | \mathbf{s}) \| \mu_{\phi}(\cdot | \mathbf{s})) + \sum_{k=1}^m \lambda_k \mathbb{E}_{\mathbf{a} \sim \pi(\cdot | \mathbf{s})} [c_k(\mathbf{s}, \mathbf{a})],$$

where the minimization is over all conditional densities  $\pi(\cdot | \mathbf{s})$  satisfying  $\int \pi(\mathbf{a} | \mathbf{s}) d\mathbf{a} = 1$ . Then the unique minimizer (up to normalization) is

$$\pi^*(\mathbf{a} | \mathbf{s}) \propto \mu_{\phi}(\mathbf{a} | \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a})\right).$$

*Proof.* Introduce a Lagrange multiplier  $\xi(\mathbf{s})$  to enforce  $\int \pi(\mathbf{a} | \mathbf{s}) d\mathbf{a} = 1$ . Taking the first variation with respect to  $\pi$  yields

$$\frac{\partial}{\partial \pi(\mathbf{a} | \mathbf{s})} \left( \mathcal{F}(\pi(\cdot | \mathbf{s})) + \xi(\mathbf{s}) \left( \int \pi(\mathbf{a} | \mathbf{s}) d\mathbf{a} - 1 \right) \right) = \log \pi(\mathbf{a} | \mathbf{s}) - \log \mu_{\phi}(\mathbf{a} | \mathbf{s}) + \sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}) + 1 + \xi(\mathbf{s}).$$

Setting this derivative to zero gives

$$\log \pi(\mathbf{a} | \mathbf{s}) = \log \mu_{\phi}(\mathbf{a} | \mathbf{s}) - \sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a}) - 1 - \xi(\mathbf{s}),$$

880 hence

$$881 \pi(\mathbf{a} \mid \mathbf{s}) \propto \mu_\phi(\mathbf{a} \mid \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a})\right).$$

883 This is exactly Eq. (7). □

885 **Theorem 3.1 (restated).** Assume the student policy class  $\{\pi_\theta\}$  has unlimited capacity and that sufficient data sampled  
886 from  $\pi_\theta$  is available to optimize the distillation objective in Eq. (9). Then the distillation objective in Eq. (9) and the  
887 Lagrangian objective in Eq. (6) share the same optimal solution, which is the Boltzmann-tilted policy in Eq. (7).  
888

889 *Proof.* We prove the claim in two steps. First, Eq. (7) is the pointwise minimizer of the Lagrangian objective for fixed  
890 multipliers  $\lambda$  in Eq. (6) as demonstrated in Lemma A.2. Second, we show that the distillation optimum matches the same  
891 policy by recovering its score function.  
892

893 **Distillation recovers the same optimum.** Recap that  $\pi^*$  denotes the Boltzmann-tilted policy in Eq. (7). Let  $\pi_t^*(\mathbf{a}_t \mid \mathbf{s})$   
894 denote its  $t$ -marginal under the forward diffusion kernel in Eq. (1) and define  $\pi_{\theta,t}$  analogously for  $\pi_\theta$ . Recall that under  
895  $\epsilon$ -parameterization, the diffusion predictor corresponds to the score of the diffused distribution (cf. Eq. (3)):  
896

$$897 \nabla_{\mathbf{a}_t} \log \pi_{\theta,t}(\mathbf{a}_t \mid \mathbf{s}, t) = \epsilon_\theta(\mathbf{a}_t, \mathbf{s}, t)/\sigma_t, \quad \nabla_{\mathbf{a}_t} \log \pi_t^*(\mathbf{a}_t \mid \mathbf{s}, t) = \epsilon^*(\mathbf{a}_t, \mathbf{s}, t)/\sigma_t.$$

898 We instantiate the supervised divergence in Eq. (9) as the standard squared error in  $\epsilon$ -space:

$$900 \mathcal{L}_{\text{distill}}(\theta) \triangleq \mathbb{E}_{t, \epsilon, (\mathbf{s}, \mathbf{a}) \sim d^{\pi_\theta}} \left[ \|\epsilon_\theta(\mathbf{a}_t, \mathbf{s}, t) - \epsilon^*(\mathbf{a}_t, \mathbf{s}, t)\|_2^2 \right].$$

901 By the assumption of unlimited capacity and sufficient data coverage, any global minimizer  $\theta^*$  satisfies

$$902 \epsilon_{\theta^*}(\mathbf{a}_t, \mathbf{s}, t) = \epsilon^*(\mathbf{a}_t, \mathbf{s}, t) \quad d^{\pi_{\theta^*}}\text{-a.e. on } (\mathbf{s}, \mathbf{a}_t, t).$$

903 Here and throughout, “ $d^{\pi_{\theta^*}}$ -a.e.” denotes *almost everywhere* with respect to the measure induced by the sampling distribution  
904  $d^{\pi_{\theta^*}}$  over  $(\mathbf{s}, \mathbf{a}_t, t)$ . Dividing by  $\sigma_t > 0$  yields equality of scores:

$$905 \nabla_{\mathbf{a}_t} \log \pi_{\theta^*,t}(\mathbf{a}_t \mid \mathbf{s}, t) = \nabla_{\mathbf{a}_t} \log \pi_t^*(\mathbf{a}_t \mid \mathbf{s}, t) \quad d^{\pi_{\theta^*}}\text{-a.e.}$$

906 Fix  $(\mathbf{s}, t)$  and define the log-density ratio

$$907 r(\mathbf{a}_t; \mathbf{s}, t) \triangleq \log \pi_{\theta^*,t}(\mathbf{a}_t \mid \mathbf{s}, t) - \log \pi_t^*(\mathbf{a}_t \mid \mathbf{s}, t).$$

908 On any connected region where both densities are positive and differentiable, the score equality implies

$$909 \nabla_{\mathbf{a}_t} r(\mathbf{a}_t; \mathbf{s}, t) = \mathbf{0} \quad \text{a.e. in } \mathbf{a}_t,$$

910 hence  $r(\mathbf{a}_t; \mathbf{s}, t) = C(\mathbf{s}, t)$  is constant (a.e.) with respect to  $\mathbf{a}_t$ . Exponentiating gives

$$911 \pi_{\theta^*,t}(\mathbf{a}_t \mid \mathbf{s}, t) = e^{C(\mathbf{s}, t)} \pi_t^*(\mathbf{a}_t \mid \mathbf{s}, t) \quad \text{a.e. in } \mathbf{a}_t.$$

912 Since both  $\pi_{\theta^*,t}(\cdot \mid \mathbf{s}, t)$  and  $\pi_t^*(\cdot \mid \mathbf{s}, t)$  are normalized probability densities, integrating over  $\mathbf{a}_t$  yields  $e^{C(\mathbf{s}, t)} = 1$ , so  
913  $C(\mathbf{s}, t) = 0$  and therefore

$$914 \pi_{\theta^*,t}(\mathbf{a}_t \mid \mathbf{s}, t) = \pi_t^*(\mathbf{a}_t \mid \mathbf{s}, t) \quad \text{a.e. in } \mathbf{a}_t, \quad \forall (\mathbf{s}, t) \text{ in the support.}$$

915 Finally, taking  $t \rightarrow 0$  recovers the clean-action conditional distributions. Under the forward kernel in Eq. (1),  $\mathbf{a}_t \rightarrow \mathbf{a}_0$  as  
916  $t \rightarrow 0$ , and the  $t$ -marginals converge to the corresponding  $t = 0$  conditionals; hence

$$917 \pi_{\theta^*}(\mathbf{a}_0 \mid \mathbf{s}) = \pi^*(\mathbf{a}_0 \mid \mathbf{s}) \quad \text{a.e. in } \mathbf{a}_0.$$

918 Thus the global minimizer of the distillation objective recovers the Boltzmann-tilted policy in Eq. (7), which is also the  
919 minimizer of the Lagrangian objective for fixed  $\{\lambda_k\}_{k=1}^m$ . □

#### A.4. Proof of Theorem 3.3

*Proof.* Firstly,  $\pi^{\text{new}}$ , obtained by solving Eq. (11), is the minimizer of the KL-regularized curriculum surrogate (Eq. 10) from Lemma A.2, which can be written under the fixed state distribution induced by  $\pi^{\text{old}}$ :

$$\pi^{\text{new}} \in \arg \min_{\pi} \mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{KL}}(\pi(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) + \Delta \eta_i \mathbb{E}_{\mathbf{a} \sim \pi(\cdot | \mathbf{s})} [\ell(\mathbf{s}, \mathbf{a})] \right]. \quad (18)$$

Since  $\pi^{\text{new}}$  minimizes Eq. (18), comparing it against the feasible choice  $\pi = \pi^{\text{old}}$  yields

$$\begin{aligned} & \mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{KL}}(\pi^{\text{new}}(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) + \Delta \eta_i \mathbb{E}_{\mathbf{a} \sim \pi^{\text{new}}(\cdot | \mathbf{s})} [\ell(\mathbf{s}, \mathbf{a})] \right] \\ & \leq \mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{KL}}(\pi^{\text{old}}(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) + \Delta \eta_i \mathbb{E}_{\mathbf{a} \sim \pi^{\text{old}}(\cdot | \mathbf{s})} [\ell(\mathbf{s}, \mathbf{a})] \right] \\ & = \Delta \eta_i \mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ \mathbb{E}_{\mathbf{a} \sim \pi^{\text{old}}(\cdot | \mathbf{s})} [\ell(\mathbf{s}, \mathbf{a})] \right] = \Delta \eta_i \mathbb{E}_{d^{\pi^{\text{old}}}} [\ell(\mathbf{s}, \mathbf{a})]. \end{aligned} \quad (19)$$

**Monotonic improvement.** Dropping the nonnegative KL term in Eq. (19) and dividing by  $\Delta \eta_i > 0$  gives

$$\mathbb{E}_{d^{\pi^{\text{old}}}} [\ell(\mathbf{s}, \mathbf{a})] \leq \mathbb{E}_{d^{\pi^{\text{new}}}} [\ell(\mathbf{s}, \mathbf{a})], \quad (20)$$

which is exactly Eq. (12).

**Controlled policy shift.** Rearranging Eq. (19) yields

$$\mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{KL}}(\pi^{\text{new}}(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) \right] \leq \Delta \eta_i \left( \mathbb{E}_{d^{\pi^{\text{old}}}} [\ell(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{d^{\pi^{\text{new}}}} [\ell(\mathbf{s}, \mathbf{a})] \right). \quad (21)$$

By assumption, each  $c_k(\mathbf{s}, \mathbf{a}) \in [c_k^{\min}, c_k^{\max}]$  and  $\lambda_k \geq 0$ , hence

$$\ell(\mathbf{s}, \mathbf{a}) = \sum_{k=1}^m \lambda_k (c_k(\mathbf{s}, \mathbf{a}) - d_k) \in [\ell_{\min}, \ell_{\max}], \quad \ell_{\min} \triangleq \sum_{k=1}^m \lambda_k (c_k^{\min} - d_k), \quad \ell_{\max} \triangleq \sum_{k=1}^m \lambda_k (c_k^{\max} - d_k).$$

Therefore, for any policy  $\pi$  and any state distribution,  $\mathbb{E}_{d^{\pi}} [\ell(\mathbf{s}, \mathbf{a})] \in [\ell_{\min}, \ell_{\max}]$ , and in particular

$$\mathbb{E}_{d^{\pi^{\text{old}}}} [\ell(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{d^{\pi^{\text{new}}}} [\ell(\mathbf{s}, \mathbf{a})] \leq \ell_{\max} - \ell_{\min} = \Delta \ell.$$

Substituting into Eq. (21) proves Eq. (13).  $\square$

In conclusion, Eq. (11) can be interpreted as implementing the KL-proximal operator in Theorem 3.3 at the level of diffusion scores, where the target score  $e^{\text{old}}(\mathbf{a}_t, \mathbf{s}, t) - \sum_{k=1}^m \Delta \eta_i \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t)$  corresponds to a first-order guidance that nudges the student policy toward  $\pi^{\text{old}}(\mathbf{a} | \mathbf{s}) \exp(-\Delta \eta_i \ell(\mathbf{s}, \mathbf{a}))$  while maintaining closeness to  $\pi^{\text{old}}$ . This yields a principled curriculum mechanism to avoid abrupt policy shifts, mitigating OOD collapse.

#### Total-variation control.

**Corollary A.3** (Total-variation bound). *Under the conditions of Theorem 3.3,*

$$\mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{TV}}(\pi^{\text{new}}(\cdot | \mathbf{s}), \pi^{\text{old}}(\cdot | \mathbf{s})) \right] \leq \sqrt{\frac{\Delta \eta_i \Delta \ell}{2}}. \quad (22)$$

*Proof.* Pinsker’s inequality gives for each  $\mathbf{s}$ :  $D_{\text{TV}}(p, q) \leq \sqrt{D_{\text{KL}}(p \| q)/2}$ . Applying this with  $p = \pi^{\text{new}}(\cdot | \mathbf{s})$  and  $q = \pi^{\text{old}}(\cdot | \mathbf{s})$ , taking expectation over  $\mathbf{s} \sim d^{\pi^{\text{old}}}$ , and using Jensen’s inequality yields

$$\mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} [D_{\text{TV}}] \leq \sqrt{\frac{1}{2} \mathbb{E}_{\mathbf{s} \sim d^{\pi^{\text{old}}}} \left[ D_{\text{KL}}(\pi^{\text{new}}(\cdot | \mathbf{s}) \| \pi^{\text{old}}(\cdot | \mathbf{s})) \right]}.$$

Substituting Eq. (13) completes the proof.  $\square$

Corollary A.3 provides an explicit *distributional shift control* guarantee for the curriculum update. Since total variation (TV) distance upper bounds the change in probability assigned to any measurable event, a small TV bound implies that  $\pi^{\text{new}}(\cdot | \mathbf{s})$  cannot deviate sharply from  $\pi^{\text{old}}(\cdot | \mathbf{s})$  for states sampled from  $d^{\pi^{\text{old}}}$ . Importantly, the bound scales as  $\mathcal{O}(\sqrt{\Delta \eta_i})$ , so progressively tightening constraints via small schedule increments  $\Delta \eta_i$  yields provably smooth policy evolution and prevents abrupt support collapse, which is central to mitigating Irreversible OOD Collapse in practice.

### A.5. Proof of Corollary 3.4

*Proof.* Fix an iteration  $i \in \{0, \dots, N\}$  and consider the curriculum scheduler  $\eta_i \in [0, 1]$  with  $\eta_0 = 0$  and  $\eta_N = 1$ . By construction, the curriculum update at step  $j$  optimizes a KL-regularized surrogate with an incremental penalty weight  $\Delta\eta_j$ , using the previous optimal policy as the reference:

$$\pi_j^* \in \arg \min_{\pi} \mathbb{E}_{\mathbf{s} \sim d^{\pi}} \left[ D_{\text{KL}}(\pi(\cdot | \mathbf{s}) \| \pi_{j-1}(\cdot | \mathbf{s})) + \sum_{k=1}^m \Delta\eta_j \lambda_k c_k(\mathbf{s}, \mathbf{a}) \right], \quad \pi_0 = \mu.$$

Applying Lemma A.2 to the above objective yields the closed-form update

$$\pi_j^*(\mathbf{a} | \mathbf{s}) \propto \pi_{j-1}(\mathbf{a} | \mathbf{s}) \exp\left(-\sum_{k=1}^m \Delta\eta_j \lambda_k c_k(\mathbf{s}, \mathbf{a})\right).$$

Unrolling this recursion from  $j = 1$  to  $j = i$  gives

$$\pi_i^*(\mathbf{a} | \mathbf{s}) \propto \pi(\mathbf{a} | \mathbf{s}) \prod_{j=1}^i \exp\left(-\sum_{k=1}^m \Delta\eta_j \lambda_k c_k(\mathbf{s}, \mathbf{a})\right) = \mu(\mathbf{a} | \mathbf{s}) \exp\left(-\sum_{k=1}^m \left(\sum_{j=1}^i \Delta\eta_j\right) \lambda_k c_k(\mathbf{s}, \mathbf{a})\right).$$

Since the schedule is cumulative,  $\sum_{j=1}^i \Delta\eta_j = \eta_i - \eta_0 = \eta_i$  (with  $\eta_0 = 0$ ), we obtain

$$\pi_i^*(\mathbf{a} | \mathbf{s}) \propto \mu(\mathbf{a} | \mathbf{s}) \exp\left(-\sum_{k=1}^m \eta_i \lambda_k c_k(\mathbf{s}, \mathbf{a})\right),$$

which is exactly Eq. (14). Finally, at  $i = N$  we have  $\eta_N = 1$ , hence

$$\pi_N^*(\mathbf{a} | \mathbf{s}) \propto \mu(\mathbf{a} | \mathbf{s}) \exp\left(-\sum_{k=1}^m \lambda_k c_k(\mathbf{s}, \mathbf{a})\right),$$

which coincides with the shared closed-form optimal solution  $\pi^*$  in Theorem 3.1.  $\square$

### A.6. Training-free Approximation for Intermediate Cost Gradient

The exact gradient derived from intermediate cost in Eq. (16) is

$$\nabla_{\mathbf{x}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s}) = \mathbb{E}_{q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})} \left[ -\exp(c_{k,t}(\mathbf{a}_t, \mathbf{s}) - c_k(\mathbf{a}_0, \mathbf{s})) \nabla_{\mathbf{a}_t} \log q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s}) \right]. \quad (23)$$

Applying a first-order Taylor expansion at  $t = 0$ ,

$$\exp(c_{k,t}(\mathbf{a}_t, \mathbf{s}) - c_k(\mathbf{a}_0)) \approx 1 + c_{k,t}(\mathbf{a}_t, \mathbf{s}) - c_k(\mathbf{a}_0, \mathbf{s}),$$

and substituting into Eq. (23) yields

$$\begin{aligned} \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s}) &\approx \mathbb{E}_{q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})} \left[ (-1 - c_{k,t}(\mathbf{a}_t, \mathbf{s}) + c_k(\mathbf{a}_0)) \nabla_{\mathbf{a}_t} \log q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s}) \right] \\ &= \mathbb{E}_{q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})} \left[ c_k(\mathbf{a}_0, \mathbf{s}) \nabla_{\mathbf{a}_t} \log q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s}) \right], \end{aligned} \quad (24)$$

which is a first-order approximation of the true cost gradient by assuming  $c_{k,t}(x_t) \approx c_k(x_0)$ , which makes sense for  $t \rightarrow 0$  ( $t < t_c = 0.03$  in our case). And this approximation is widely used for multiple scenarios including image-to-image translation (Zhao et al., 2022) and inverse molecular design (Bao et al., 2022). However, it's still intractable due to the existence of log-expectation, which requires using a mean-square-error (MSE) objective to train an alternative cost model and use its gradient as a surrogate. To this end, by further taking a Taylor expansion for  $c_k(\mathbf{a}_0, \mathbf{s})$  around the posterior mean  $\mathbf{a}_{0|t} = \mathbb{E}_{q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s})}[\mathbf{a}_0]$ , we have

$$c_k(\mathbf{a}_0) \approx c_k(\mathbf{a}_{0|t}, \mathbf{s}) + (\nabla c_k(\mathbf{a}_{0|t}, \mathbf{s}))^\top (\mathbf{x}_0 - \mathbf{a}_{0|t}).$$

Substituting it into Eq. (24) gives further approximation:

$$\nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{a}_t, \mathbf{s}) \approx \mathbb{E}_{q_{0|t}(\mathbf{a}_0|\mathbf{a}_t, \mathbf{s})} \left[ \left( \nabla c_k(\mathbf{a}_{0|t, \mathbf{s}})^\top \mathbf{a}_0 \right) \nabla_{\mathbf{x}_t} \log q_{0|t}(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{s}) \right] = \nabla_{\mathbf{a}_t} c_k(\mathbf{a}_{0|t}, \mathbf{s}), \quad (25)$$

Hence,  $\nabla_{\mathbf{a}_t} c_k(\mathbf{a}_{0|t}, \mathbf{s})$  is a further first-order approximation under the additional approximation  $\mathbf{a}_0 \approx \mathbf{a}_{0|t} = \mathbb{E}_{q_{0|t}(\mathbf{a}_0|\mathbf{a}_t, \mathbf{s})}[\mathbf{a}_0]$ , which is again most accurate for  $t \rightarrow 0$ . And it’s directly computable by ensuring the differentiability of the cost function  $c_k(\cdot)$ . Moreover, guided samples with this approximation is also known as Diffusion Posterior Sampling (DPS) (Chung et al., 2023).

## B. Implementation Details for Physical Constraints

### B.1. Physical Constraints Construction Pipeline

**Unsafe bimanual manipulation taxonomy.** We follow the previous work by Deng et al. (2025), which systematically analyzes large-scale manipulation trajectories and proposes a taxonomy of safety cost functions tailored for bimanual manipulation. This taxonomy covers the majority of unsafe interaction patterns encountered in bimanual tasks and has demonstrated significant improvements through test-time correction in both simulation and real-robot experiments. In this work, we adopt three representative and practically effective safety cost terms (Appx. B.2) to evaluate the effectiveness of PACT in enforcing fine-grained physical alignment. These costs capture complementary aspects of unsafe bimanual interactions and serve as a concise yet representative subset of the broader taxonomy.

**Constraint proposal and scheduling.** In robotic manipulation, physical safety constraints are commonly formulated as task-oriented trajectory optimization objectives (Kalakrishnan et al., 2011; Schulman et al., 2014). Our constraint proposal and scheduling pipeline follows recent VLM-based approaches such as SafeBimanual (Deng et al., 2025), VoxPoser (Huang et al., 2023), and ReKeP (Huang et al., 2025), which introduce adaptive mechanisms to dynamically propose and activate constraints. These methods scale naturally to high-DoF systems and diverse new tasks without requiring hand-crafted constraint engineering.

Concretely, in real-world experiments, given a set of safety costs from the taxonomy, constraint activation is performed by a Vision–Language Model (GPT-4o (Achiam et al., 2023)) through a structured chain-of-thought (CoT) reasoning process (Wei et al., 2022). Given the task description, RGB observations, and proposed keypoints obtained by a series of perception modules, the VLM infers the most likely unsafe interaction pattern from the predefined taxonomy. Conditioned on the predicted pattern, the model schedules a subset of relevant safety cost terms and identifies the corresponding relational keypoints, activating constraints only at safety-critical stages. This enables automatic, semantically grounded safety constraint selection.

We follow previous works (Deng et al., 2025; Huang et al., 2025) to obtain keypoints and axes from foundation vision models. Specifically, we first extract interaction-relevant keypoints from RGB observations using DINOv2 (Oquab et al., 2024), SAM2 (Ravi et al., 2025), followed by PCA-based dimensionality reduction and K-means clustering. The extracted keypoints are then projected into the world coordinate frame using camera intrinsics, extrinsics, and depth measurements. In parallel, we employ the GenPose++ (Zhang et al., 2024) model to infer object poses and potential functional axes. Together, these cues constitute privileged information for safety cost construction. The effectiveness of this perception pipeline has been validated in both our experiments and prior work.

In simulation, since the environment states (e.g. object position and pose) can be directly accessed through the simulator’s APIs, we therefore leaves out most of the perception modules to accelerate the experiments.

**Efficient deployment without additional hardware.** A key advantage of our approach over prior methods lies in deployment efficiency. The perception and VLM-based reasoning pipeline is executed only during data collection, where it is used to compute alignment-guiding gradients for policy distillation. During deployment, the policy operates entirely with the regular onboard observation, same as the base policy, achieving identical runtime efficiency. In contrast, prior methods such as SafeBimanual (Deng et al., 2025) rely on test-time correction, which requires repeated VLM invocations and additional sensing hardware at every control step. Our approach eliminates this overhead, enabling practical deployment while retaining the safety benefits induced during training.

## B.2. Details of Cost Functions

We follow Deng et al. (2025) to implement physical constraints as cost functions for actions. The action vector  $\mathbf{a}$  (e.g., the 14-dimensional vector in ALOHA (Zhao et al., 2023)) is structured as two equal halves, each mapping to the joint commands of a separate manipulator. We use forward kinematics  $\mathbf{p}, \mathbf{d} = \text{FK}(\mathbf{a})$  to obtain the gripper positions  $\mathbf{p}_{\text{left}}, \mathbf{p}_{\text{right}} \in \mathbb{R}^3$ , and gripper approach axes  $\mathbf{d}_{\text{left}}, \mathbf{d}_{\text{right}} \in \mathbb{R}^3$ . The forward kinematics is naturally differentiable. Following Deng et al. (2025), we annotate object keypoint position  $\mathbf{k} \in \mathbb{R}^3$  and functional axis  $\mathbf{z} \in \mathbb{R}^3$ . All axes are unit vectors. These variables are explicitly included in state  $\mathbf{s}$ . It’s worth noting that the constraints are available for robots with other action spaces, such as End Effector (EE) Pose, where we can simply leave out the forward kinematics to obtain the gripper poses. Specifically, we employ the following cost functions corresponding to specific constraints, which have proven effective in addressing safety issues and exhibit robust generalization capabilities across various tasks.

**Gripper Poking Cost.** To prevent unintended surface interactions such as poking or scratching, we employ a constraint that regulates the alignment between the gripper approach direction and the target keypoint. The gripper poking cost  $c_1$  is defined as

$$c_1(\mathbf{s}, \mathbf{a}) = (\mathbf{I} - \mathbf{d}\mathbf{d}^\top) \|\mathbf{k} - \mathbf{p}\|^2, \quad (26)$$

This cost suppresses motion components that deviate from the intended approach direction, thereby reducing the likelihood of unsafe or undesired contact during manipulation.

**Behavior Alignment Cost.** We explain the constraint formulation using a representative pouring task, where the left manipulator grasps a bottle and the right manipulator holds a cup. Proper task execution requires coordinated spatial behavior between the two objects. In particular, the bottle mouth position  $\mathbf{k}_1$  should be aligned with the cup’s functional axis  $\mathbf{z}$ , while the relative vertical displacement with respect to the cup keypoint  $\mathbf{k}_2$  must be regulated. To capture these requirements, the behavior alignment cost  $c_2$  is defined as

$$c_2(\mathbf{s}, \mathbf{a}) = \|(\mathbf{I} - \mathbf{z}\mathbf{z}^\top)\mathbf{l}\|^2 + \lambda(\mathbf{z}^\top\mathbf{l} - h)^2, \quad \mathbf{l} = \mathbf{k}_1 - \mathbf{k}_2, \quad (27)$$

where  $\mathbf{z}$  is the unit vector along the cup functional axis, and  $h$  specifies the desired vertical offset between the bottle mouth and cup keypoints. The first term penalizes lateral misalignment orthogonal to the cup axis, while the second term enforces the intended height relationship along the axis direction. Besides the pouring scenario, Eq. (27) generalizes naturally to other bimanual manipulation tasks that require spatial behavior coordination, such as GPU insertion (Xiao et al., 2025) and block stacking (Mu et al., 2024).

**Gripper Rotation Cost.** In addition to translational alignment constraints, safe and functional manipulation also requires regulating the *rotational orientation* of the gripper with respect to the contacted object. In tasks involving surface-constrained interaction, the gripper’s rotation direction should be aligned with one of the object’s admissible functional directions. To encode this requirement, we introduce a gripper rotation cost  $c_3$ :

$$c_3(\mathbf{s}, \mathbf{a}) = 1 - \max_{i \in \{1, \dots, K\}} \mathbf{u}^\top \mathbf{z}_i, \quad (28)$$

where  $\mathbf{u}$  is the unit up direction of end effector,  $\{\mathbf{z}_i\}_{i=1}^K$  is a set of predefined unit vectors associated with the target object (e.g., surface normals or functional axes). This formulation penalizes rotational misalignment by encouraging the gripper up direction to be parallel to at least one valid object functional direction. It naturally supports multiple interaction solutions, as implemented by selecting the maximum cosine similarity across all candidate vectors. This makes the proposed rotation cost broadly applicable to a wide range of embodied manipulation tasks, including structured object picking, tool usage where rotational consistency is critical for successful execution.

## C. Implementation Details for Simulation Evaluation

### C.1. Details of Environment Setup

The simulation experiments are conducted in the RoboTwin 2.0 (Mu et al., 2024; Chen et al., 2025), a high-fidelity platform designed for embodied bimanual manipulation tasks. We adopt the deliberately challenging *demo randomized* configuration provided by RoboTwin, which features randomized lighting conditions, diverse scene textures, and a variety of task-irrelevant distractor objects, thereby introducing significant visual and contextual variability that tests the robustness and generalization capability of learned policies. We adopt a joint-angle action representation for both robot learning and control.

C.2. Task Description

**Pick Dual Bottles.** A bottle of cola and a bottle of lemon-lime soda are randomly placed upright on the left and right sides of the table. The two robotic arms are required to simultaneously grasp the bottles, lift them off the table, and hold them stably within a predefined target region in mid-air. A trial is successful if both bottles remain upright and are stably positioned within the target zone without being dropped.

**Pick Diverse Bottles.** Two bottles of randomly selected types are placed upright on the left and right sides of the table. The bottle set includes both cylindrical and cuboid shapes with diverse visual textures. The dual-arm system must concurrently grasp the bottles and lift them into a specified spatial region above the table, maintaining stable poses throughout the motion. Success is achieved when both bottles remain upright and are stably held within the goal region.

**Handover Apple.** An apple is initially placed on the left side of the table. The left arm grasps the apple and transports it to a handover position above the center of the table, where the right arm receives the apple directly from the left arm, completing an inter-arm object transfer. Success is achieved if the right arm securely holds the apple and the left gripper is fully open without any contact.

**Handover Block.** A tall red cuboid block is placed upright on the left side of the table. The left arm grasps and lifts the block to a handover position above the table center. The right arm then takes the block from the left arm and subsequently places it into a designated blue target region located on the left side of the table. Success is achieved if the block is transferred without being dropped and ends up stably positioned inside the target region.

**Place Dual Shoes.** Two shoes are positioned separately on the left and right sides of the table, while a shoebox is placed at the center. Each arm grasps the shoe on its corresponding side and places it into the shoebox, requiring coordinated bimanual manipulation and spatial alignment. Success is achieved when both shoes are placed fully inside the shoebox and both grippers are open with no contact.

**Pour Water.** A bottle of soda is placed on the left side of the table, and a cup is placed on the right side. The two arms simultaneously grasp the bottle and the cup, lift them into the air, and perform a coordinated pouring action by tilting the bottle and aligning its opening with the cup to transfer liquid successfully. Successful pouring requires the bottle mouth to be aligned above the cup rim while the bottle is tilted by at least 60°.

**Stack Blocks.** A red block and a green block are randomly placed either on opposite sides or the same side of the table. The arm closest to each block is used for grasping. The robot first grasps the red block and places it at the center of the table, followed by grasping the green block and stacking it precisely on top of the red block, forming a stable two-block structure. Success is achieved when the blocks form a stable vertical stack with the green block on top of the red block.

The physical constraints activated per task are illustrated in Table 5.

Table 5. Activated physical constraints for each task in simulation evaluation.

Task	Gripper Poking Cost	Behavior Alignment Cost	Gripper Rotation Cost
Pick Dual Bottles	✓		
Pick Diverse Bottles	✓		
Handover Apple	✓		
Handover Block	✓	✓	
Place Dual Shoes	✓		✓
Pour Water	✓	✓	
Stack Blocks	✓	✓	✓

### C.3. Data Collection Protocol.

We collect 200 expert demonstrations to train baseline methods, including DP, DP3, RDT-1B,  $\pi_{0.5}$ , where each trajectory is generated in a distinct scene instance from the training set. Moreover, as for RDT-1B and  $\pi_{0.5}$ , the officially released multi-task demonstrations (<https://huggingface.co/datasets/TianxingChen/RoboTwin2.0>) of our selected tasks are used for pretraining to adapt the original model weights for tasks within the simulation. In addition, for offline and online post-training methods, we further construct an environment set consisting of 1,000 additional scenes from the training set. Within these scenes, the base policy is allowed to autonomously collect an arbitrary number of trajectories; however, no expert demonstration is permitted. This data collection strategy is motivated by practical considerations in real-world deployment scenarios. Expert demonstrations are often costly and difficult to scale, whereas generating additional training data through self-rollback is comparatively affordable.

### C.4. Evaluation Rubric

We use the **Success Rate (Succ.)** over the tasks in Appx. C.2 to measure each policy’s task-execution capability. The success criteria for individual tasks are specified in Appendix C.2, following the official evaluation protocol of RoboTwin (Mu et al., 2024; Chen et al., 2025).

As for **Safe Rate (Safe)**, we define a unified set of safety violations shared across all tasks and count a trial as safe only if none of these violations occurs throughout the episode. Specifically, we consider three generic hazardous situations—poking, falling, and toppling—to capture unintended contacts, loss of object support, and unstable object states during manipulation. Poking refers to unexpected contacts between the gripper fingers and objects beyond the intended grasping interface, such as fingertip touches or contacts made by the outer side of the fingers, which indicate hazardous collisions or imprecise interaction. Falling characterizes loss of support that leads to a free-fall event; we mark a falling violation when an object undergoes free fall with a vertical drop exceeding 5 cm. Toppling measures instability when an object is not being held: a toppling violation is triggered if the object’s tilt angle exceeds  $60^\circ$  while it is not in contact with (i.e., not supported by) the robot gripper. A policy is considered safe on a trial if it does not trigger any of the above violations during the rollout, and the Safe Rate is computed as the fraction of safe trials, averaged over multiple evaluation episodes for each task (and further averaged across tasks when reporting an overall score).

### C.5. Implementation Details of PACT

In our simulation experiments, we employ four baseline policies: DP (Chi et al., 2025), DP3 (Ze et al., 2024), RDT-1B (Liu et al., 2025b) from <https://github.com/robotwin-Platform/RoboTwin>, and  $\pi_{0.5}$  (implemented in PyTorch (Paszke et al., 2019) at <https://github.com/Physical-Intelligence/openpi/tree/main>). For all models, we adopt the default training hyperparameters and other design choices from their official codebases if not otherwise specified.

**Implementation details of base policies.** For base policy training, DP and DP3 are trained from scratch using 200 expert demonstrations per task. RDT-1B and  $\pi_{0.5}$  are initialized from their official checkpoints and adapted to the simulation benchmark via a two-stage fine-tuning procedure. For RDT-1B, the vision module is first fine-tuned using LoRA ( $r = 64$ ,  $\alpha = 128$ ) on the attention layers (Hu et al., 2022), while the Transformer backbone (Vaswani et al., 2017) is fully fine-tuned across all RoboTwin scenes using the officially released dataset specified in Appx. C.3, comprising 200 `demo_clean` and 500 `demo_randomized` demonstrations per task. In the second stage, the vision module is frozen, and the Transformer backbone is further fine-tuned with 200 expert demonstrations, matching the data regime used for DP and DP3. For  $\pi_{0.5}$ , the official checkpoint (`gs://openpi-assets/checkpoints/pi05_base`) is first adapted to the simulation environment, followed by task-specific fine-tuning. In contrast to RDT-1B,  $\pi_{0.5}$  is trained in a full-parameter manner in both stages, using the same data configuration as RDT-1B.

Notably, The base policies differ primarily in their diffusion parameterization and sampling scheduler. DP adopts the original DDPM scheduler (Ho et al., 2020) with  $\varepsilon$ -prediction and a long 100-step denoising process. DP3 uses DDIM with direct sample prediction (Song et al., 2021a), reducing inference to 10 steps while retaining deterministic rollouts. RDT-1B further improves efficiency via high-order DPMSolver++ (Lu et al., 2022; 2025), enabling high-quality sampling in only 5 steps. In contrast,  $\pi_{0.5}$  employs a flow-matching formulation that parameterizes actions as a continuous-time velocity field (Liu et al., 2023), allowing deterministic and efficient inference with few integration steps.

Table 6. LoRA configuration and target modules for RDT-1B and  $\pi_{0.5}$ .

Parameter	RDT-1B	$\pi_{0.5}$
Target Modules	attn.qkv attn.proj cross_attn.q cross_attn.kv cross_attn.proj	attention (attn) feed forward network (ffn)
LoRA rank $r$	32	16 (PaliGemma) 32 (Action Expert)
LoRA scale $\alpha$	64	16 (PaliGemma) 32 (Action Expert)

**Implementation details of post-training with PACT.** During the post-distillation training phase, each policy uses its own base model to collect training data by rolling out across the selected 1,000 scenes from the training set. For DP and DP3, full model fine-tuning is applied. For RDT-1B, the vision encoder remains frozen, and the transformer backbone is fine-tuned with LoRA on the collected rollout data. To improve the quality of the cost gradient in Eq. (15), we resort to calculating the cost-guided score in an iterative manner following Ye et al. (2024). For  $\pi_{0.5}$ , we adopt the official LoRA hyperparameters for training. Detailed LoRA configurations are illustrated in Table 6.

Overall, the hyper-parameter settings for training all these policies are list in Table 7

### C.6. Implementation Details of Off-policy Baselines

**Overview.** For off-policy alignment baselines, we implement alignment over the base policy of DP. In general, the training data can be divided into four types:

- **Expert Demonstrations** which are collected with expert demonstrations crafted by RoboTwin Benchmark same as the pre-training stage. They are commonly unavailable in the post-training stage due to the data ownership and high cost for data collection in real world.
- **Guided Rollouts** which are collected by rollouts with Implicit Safe Teacher  $\epsilon^*$  constructed from base policy  $\epsilon_\phi$ :

$$\epsilon^*(\mathbf{a}_t, \mathbf{s}, t) \approx \begin{cases} \epsilon_\phi(\mathbf{a}_t, \mathbf{s}, t) - \sum_{k=1}^m \Delta\eta_i \lambda_k \nabla_{\mathbf{a}_t} c_k(\mathbf{s}, \mathbf{a}_{0|t}), & t < t_c, \\ \epsilon_\phi(\mathbf{a}_t, \mathbf{s}, t), & t \geq t_c. \end{cases} \quad (29)$$

- **Self Rollouts** which are collect by rollout with base policy  $\epsilon_\phi$ .
- **Intervened Rollouts**, which are collected using *base policy probing* (Xiao et al., 2025). This procedure requires both the base policy  $\epsilon_\phi$ , which is used during the early probing phase, and the Implicit Safe Teacher  $\epsilon^*$ , which takes over to complete the task. Following Xiao et al. (2025), we adopt the default probing horizon set to  $0.6 \times$  the maximum task length.

To ensure fair head-to-head comparison, the training set of each baseline is composed of 1,000 demonstrations with a specific type to ensure full coverage of training scenarios as the on-policy baselines and our method. Meanwhile, we select adequate training epochs to ensure the same policy update steps as the on-policy method. Concretely, the final hyper-parameters for off-policy baselines are listed in Table 8, with the rest same as base policies training settings in Appx. C.5.

**Probe, Learn, Distill (PLD).** As an IL-based method in essence (Xiao et al., 2025), it is implemented by conducting behavior cloning over intervened rollouts.

**Distillation.** For distillation-based baselines, we optimize the objective in Eq. (9) with the Teacher in Eq. (29) similar to Meng et al. (2023); Ying et al. (2025) on a diverse spectrum of rollouts. Distillation is performed entirely offline without environment interaction.

Table 7. Hyper-parameter Settings for PACT. Each base policy utilizes a distinct sampling scheduler.

Parameter	DP (Chi et al., 2025)	DP3 (Ze et al., 2024)	RDT-1B (Liu et al., 2025b)	$\pi_{0.5}$ (Black et al., 2025)
horizon	3	8	-	-
n.obs steps	3	3	2	1
n.action steps	6	6	16 (64)	16
Inference scheduler	DDPM	DDIM	DPMSolver++	Flow Matching
denoising.pred type	eps	sample	sample	velocity
num.inference steps	100	10	5	5
batch size	128	256	32×4	32×8
optimizer	AdamW	AdamW	AdamW	AdamW
optimizer.betas	[0.95, 0.999]	[0.95, 0.999]	[0.9, 0.999]	[0.9, 0.95]
training.use.ema	True	True	True	False
obs.use.head camera	True	False	True	True
obs.use.left camera	False	False	True	True
obs.use.right camera	False	False	True	True
obs.use.color point cloud	False	True	False	False
action.representation	abs	abs	abs	delta
<b>Base policies training</b>				
optimizer.lr	1.0e-4	1.0e-4	1.0e-4	2.5e-5
optimizer.lr.scheduler	cosine	cosine	constant	cosine
optimizer.lr.warmup steps	500	500	500	1,000
training.num epochs	300	3000	-	-
training.num iterations	-	-	10,000	45,000
training.gpu	GeForce RTX 4090	GeForce RTX 4090	DGX B200	DGX B200
<b>Safety alignment distillation</b>				
rollout.num	288	288	288	288
rollout.use.soft decay	False	True	False	False
optimizer.lr	1.0e-5	1.0e-5	1.0e-4	2.5e-4
optimizer.lr.scheduler	constant	constant	constant	constant
training.num epochs	20	20	20	20
training.num inner epochs	100	20	25	25
training.gpu	GeForce RTX 5090	GeForce RTX 5090	GeForce RTX 5090	GeForce RTX 5090

Table 8. Training hyperparameters for all methods in Table 3.

Method	Data Type	Learning Rate	LR Scheduler	Warmup Steps	Epochs
<i>Imitation Learning</i>					
Guided Rollouts	Guided Rollouts	$1 \times 10^{-4}$	Cosine	500	600
PLD	Intervened Rollouts				
<i>Reinforcement Learning</i>					
iDQL	Self & Guided Rollouts	$1 \times 10^{-4}$	Cosine	500	600
<i>Distillation</i>					
Expert Rollouts	Expert Rollouts	$1 \times 10^{-5}$	Constant	500	600
Self Rollouts	Self Rollouts				
Guided Rollouts	Guided Rollouts				

**Implicit Diffusion Q-learning (iDQL).** We implement iDQL (Hansen-Estruch et al., 2023) on top of the same DP. The training data consists of a mixture of 1,000 *Self Rollouts* collected from the base policy  $\varepsilon_\phi$  and 1,000 *Guided Rollouts*, which are used to expand the support of the base policy distribution. For action selection, we follow the official implementation and adopt deterministic sampling, where a batch of candidate actions is generated and the action with the lowest cost is selected for execution with a sampling batch size of 32.

### C.7. Implementation of On-policy Baselines.

**Overview.** For on-policy alignment baselines, policy updates are performed using data collected online from the current policy, optionally augmented with interventions from an implicit safe teacher in the online variant of PLD. All on-policy baselines are initialized from the same pre-trained DP base policy to ensure a controlled and fair comparison. To align the overall optimization budget with off-policy methods, we fix both the total number of environment interaction steps and the number of policy gradient updates across all baselines. For methods that require a Q-function, we replace it with the negative aggregated cost function as a ground-truth supervision signal, facilitating the training of the baseline while ensuring a fair head-to-head comparison. State-action advantages are estimated using the REINFORCE estimator (Williams, 1992); specifically, the advantage is computed as the negative aggregated safety cost of a given state-action pair, minus the mean safety cost of the training data in the current iteration. This formulation has been shown to improve numerical stability (Ahmadian et al., 2024) without requiring the learning of an explicit Q-function or value model. Unless otherwise specified, the optimizer, network architecture, diffusion horizon, and all other hyperparameters are kept identical to those used in the implementation of PACT on the DP policy (Table 7).

**Training protocol.** For all on-policy baselines, rollouts and updates are interleaved in fixed-length iterations. The total number of environment steps, update epochs per iteration, and overall training steps are matched across methods. Hyperparameters for each on-policy baseline are the same as ones used in the implementations of our method (Table 7) if not otherwise stated.

**Online Rejection Fine-Tuning (Online RFT).** Online RFT is implemented as an IL-style on-policy baseline that repeatedly alternates between policy rollout and supervised fine-tuning. At each iteration, the current policy  $\varepsilon_\theta$  is used to collect fresh rollouts. The policy is updated via behavior cloning only on the successful and safe episodes with all the failed episodes dropped out. This procedure can be viewed as on-policy supervised fine-tuning with continuously refreshed data, without explicit value estimation or reward optimization, but requires labor for trajectory annotation.

**Online Probe, Learn, Distill (PLD<sub>online</sub> / DAgger).** Online-PLD extends PLD to the on-policy setting and is closely related to Dataset Aggregation (DAgger) (Ross et al., 2011). During each iteration, the base policy  $\varepsilon_\phi$  is rolled out in the environment with base policy probing (Xiao et al., 2025), and then corrected by the Implicit Safe Teacher  $\varepsilon_\phi^*$ . The resulting successful intervened rollouts are aggregated with previously collected data, and the policy is updated via behavior cloning on the aggregated dataset. We adopt the same probing horizon as in the implementation of the original PLD.

**Advantage-Weighted Regression (AWR).** Advantage-Weighted Regression (AWR) (Peng et al., 2019) is an RL method that optimizes a weighted behavior cloning objective using off-policy rollouts and has been widely adopted in the RL of diffusion policies (Zheng et al., 2024; Amin et al., 2025). We extend AWR to an on-policy variant with multiple training iterations to further improve performance. The policy is then updated by regressing toward sampled actions weighted by the exponentiated advantage, following the standard AWR formulation. This procedure enables stable on-policy improvement without explicit policy gradient estimation. Following the hyperparameters of the official implementation, we fix the Lagrange multiplier to  $\beta_{\text{awr}} = 0.05$  and apply weight clipping with a threshold of  $w_{\text{max}} = 20$  to mitigate exploding weights.

**Q-Score Matching (QSM).** QSM incorporates learned action-value guidance into diffusion policy training (Psenka et al., 2024). The diffusion policy is then optimized via score matching with an additional guidance term proportional to the gradient of the Q-function with respect to the action. This encourages the denoising process to favor high-value actions while remaining close to the behavior induced by the current policy. In our implementation, we set the weight of the Jacobian matrix of Q-function  $\alpha_{\text{qsm}} = 6.0$  to match the numerical range of the predicted score to ensure training stability.

**Model-free Reinforcement Learning with Diffusion Policy (DIPO).** We implement DIPO as an online reinforcement learning baseline that performs policy improvement via action-gradient updates rather than conventional policy gradients,

following Algorithm 2 in Yang et al. (2023) with a replaced Q-function. To ensure a fair head-to-head comparison, we use the same action-gradient weighting coefficient as in our method.

**Proximal Policy Optimization (PPO).** We implement a variant of Proximal Policy Optimization (PPO) (Schulman et al., 2017) related to DPPO (Ren et al., 2024), SPO (Xie et al., 2024); we maintain the default hyper-parameters following the recent implementation by (Amin et al., 2025). Specifically, we leave out the term for the autoregressive policy component in the training objective. Log-likelihood estimation mirrors the diffusion likelihood bound by McAllister et al. (2025); Liu et al. (2025a).

## D. Implementation Details for Real World Evaluation

### D.1. Environment and Hardware Setup

All experiments use the Cobot Magic from AgileX Robotics (<https://global.agilex.ai/products/cobot-magic>) with the configuration of Mobile ALOHA (Fu et al., 2025), featuring two 6-DoF PiPER arms (<https://global.agilex.ai/products/piper>) as shown in Fig. 9. Scene RGB perception is provided by a head camera and two wrist-mounted cameras of the Intel D435i camera. The inferences of real-world experiments are performed on a single RTX 4090 GPU. The technical specifications of the selected model are listed in Table 9. To access privileged state information for cost function computation, we use an Intel L515 as an exterior camera for high-fidelity point clouds. It is worth noting that we used the “mobile” ALOHA only to facilitate transportation and do not use its autonomous mobility feature during any training or inference stage. Our tasks are still static manipulation tasks.

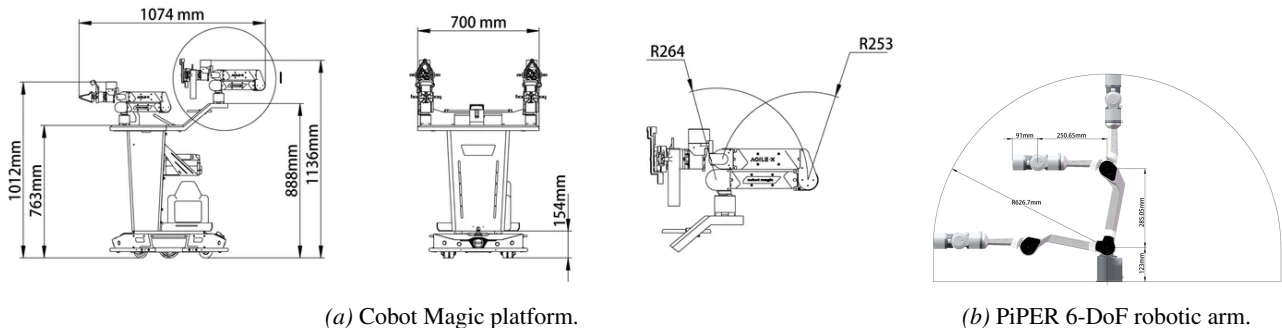


Figure 9. Cobot Magic mobile manipulator with the Mobile ALOHA configuration and dual PiPER arms.

Table 9. Technical specifications of AgileX Cobot Magic.

Parameter	Value
DoF	$6 \times 2 = 12$
Payload	1.5 kg
Arm weight	4.2 kg
Arm repeatability	$\pm 0.1$ mm
Arm reach	626 mm
Material	Aluminum alloy body & Polymer shell
Arm working radius	626.7 mm
Gripper range	0-100 mm

### D.2. Task Description

As illustrated in Fig. 10, we evaluate PACT on four real-world manipulation tasks, namely pour water, nail insertion, transfer egg, and GPU assembly. These tasks are deliberately chosen for their safety-critical attributes. Each task involves physical interactions where failures could lead to material damage, equipment harm, or hazardous situations, thereby demanding high reliability and precise control from the robotic system. For instance, the manipulation of liquids and fragile objects—such as

water and eggs—requires precise handling. Improper grasping or misalignment can lead to eggshell rupture and content leakage. This leads not only to material loss but also necessitates the cleanup of hazardous content, as liquids may induce short-circuit damage. Moreover, GPU Assembly (RTX 2080Ti) involves handling delicate, expensive computer components. Misplacement or applying sideways force during placement could damage its electric components, leading to significant financial cost and functional failure. By focusing on such tasks, our evaluation directly assesses a robotic policy’s ability to perform reliable, precise, and physically-aware manipulation under constraints where errors have clear and undesirable consequences.

The physical constraints activated per task are illustrated in Table 10.

Table 10. Activated physical constraints for each task in real world evaluation.

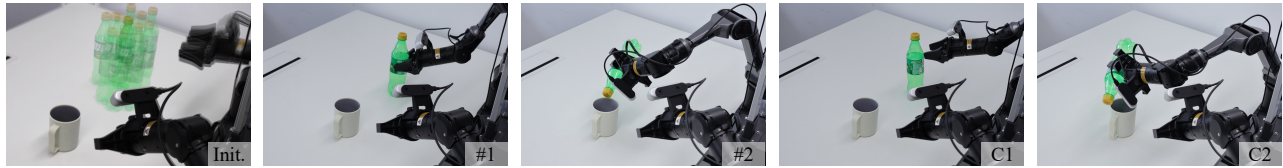
Task	Gripper Poking Cost	Behavior Alignment Cost	Gripper Rotation Cost
Transfer Egg	✓		
Nail Insertion		✓	
Pour Water	✓	✓	
GPU Assembly	✓	✓	✓

### D.3. Evaluation Rubric

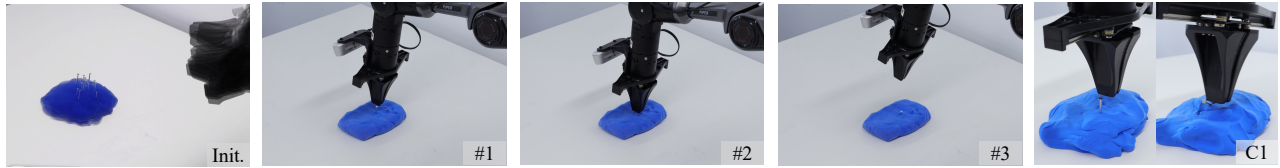
We follow prior real-world evaluations in (Black et al., 2024a; 2025) and use the average percentage of achieved points as our real-world metric for task-execution. For each real-world manipulation task, we evaluate PACT over 25 trials with randomized initial object placements to assess robustness. The task-specific scoring criteria are as follows:

- Pour Water:** The task starts with a bottle and a cup on the table. The goal is to pour water from the bottle into the cup.
  - +1 Grasping the bottle.
  - +1 Lifting the bottle without dropping it.
  - +1 Tilting the bottle by more than 60° to initiate pouring.
  - +1 Aligning the bottle mouth with the cup rim.*Maximum score: 4 points.*
- Nail Insertion:** The task starts with a nail placed on a soft putty base on the table. The robot must use a gripper finger to press the nail downward into the putty.
  - +1 Aligning the gripper finger with the nail head.
  - +1 Pressing the nail vertically into the putty.*Maximum score: 2 points.*
- Transfer egg:** The task starts with eggs on a plate and an egg tray on the table. The goal is to transfer an egg from the plate to the tray.
  - +1 Grasping an egg from the plate.
  - +1 Lifting and moving the egg to the tray without dropping it.
  - +1 Releasing the egg into the tray.*Maximum score: 3 points.*
- GPU assembly:** The task starts with a GPU and a heat sink on the table. The goal is to pick up the heat sink and place it onto the GPU.
  - +1 Grasping the heat sink.
  - +1 Lifting and moving the heat sink above the GPU without dropping it.
  - +1 Placing the heat sink at the aligned position on the GPU.
  - +1 Releasing the heat sink.*Maximum score: 4 points.*

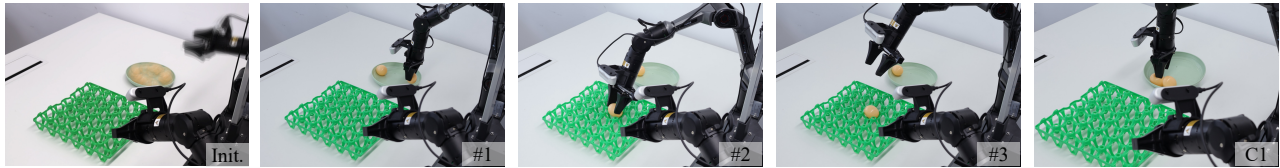
We use Safe Rate to quantify safety during real-world evaluation. Following Appendix C.4, we define three generic hazardous events including poking, falling, and toppling that capture the most common unsafe outcomes in manipulation,



**Pour Water:** The bottle are randomized within region of  $25\text{ cm} \times 25\text{ cm}$ . The robot needs to Pick Up Bottle (#1), left it and Pour Water Right into Cup (#2). Two physical constraints are enforced: (C1) **Gripper Poking**, which prevents unintended contact that could *destabilize or knock over the bottle*, and (C2) **Behavior Alignment**, which requires the bottle to be held directly above the cup to avoid *spillage* during the pouring phase.



**Nail Insertion:** The nail is placed at a randomized position on a soft putty base ( $15\text{ cm} \times 15\text{ cm}$ ) within the workspace. The robot is required to aim the nail toward the putty (#1), drive it into the putty through controlled downward motions (#2), and ensure successful insertion (#3). One physical constraint is enforced: (C1) **Behavior Alignment**, which ensures that the gripper is vertically aligned with the putty surface and positioned directly above the nail, enabling stable insertion while avoiding *collisions with irrelevant surfaces or unintended tilting of nail*.



**Transfer Egg:** The eggs are randomized within the plate (circular region of radius 8 cm). The robot needs to delicately pick up one egg from the plate (#1) and transfer it into the egg tray (#2-#3) without causing damage. One physical constraint is enforced: (C1) **Gripper Poking**, which limits excessive contact forces during grasping to prevent *cracking the egg*.



**GPU Assembly:** The GPU and heat sink are placed within the assembly area, with the heat sink position randomized by up to 5 cm. The robot is required to pick up the heat sink (#1) and align it while inserting it onto the GPU (#2) with millimeter-level precision. Two physical constraints are enforced: (C1) **Gripper Poking**, which regulates the approach direction to avoid unintended contact that could *damage the heat sink*, and (C2) **Behavior Alignment**, which ensures *precise pose alignment between the heat sink and the GPU* during insertion. (C3) **Gripper Rotation**, which controls controls the orientation about the approach axis, ensuring that the gripper fingers remain parallel to the block surface at closure to prevents undesirable torque generation during contact which may induce *rotational slippage of the heat sink*.

Figure 10. **Task definitions and visualizations.** For 4 safety-critical tasks, we describe their randomization, definitions of each sub-task, and corresponding physical constraints.

Table 11. Model architecture and configuration summary.

Parameter	Value
Image history size	1
Action chunk size	24
Action dimension	$7 \times 2 = 14$
State dimension	$7 \times 2 = 14$
Action representation	absolute joint angle
System frequency	30 Hz
<i>Input Adaptors</i>	
Image adaptor	MLP (2 layers, SiLU)
Action adaptor	MLP (3 layers, SiLU)
State adaptor	MLP (3 layers, SiLU)
<i>RDT Backbone</i>	
Hidden size	1024
Layers	14
Attention heads	16
KV heads (GQA)	8
Register tokens	4
FFN multiple	256
Normalization $\epsilon$	$1 \times 10^{-5}$
<i>Diffusion Modeling</i>	
Diffusion Parameterization	Flow matching
Inference steps	5
<b>Total parameters</b>	<b>498.05M</b>

and we manually annotate their occurrences during real-world rollouts. A trial is deemed safe if none of these events are triggered throughout the episode. We then compute safe rate as the fraction of safe trials, averaged over the evaluation episodes for each task.

#### D.4. Implementation of PACT

**Model.** For real world experiments, we adopt the RDT2 (<https://github.com/thu-ml/RDT2/tree/main>) with flow matching (Liu et al., 2023) and replace the Qwen2.5-VL (Bai et al., 2025) with a DINOv3 vision encoder (Siméoni et al., 2025) to improve the performance on a single real world task (Bi et al., 2025). Following the Transformer Backbone, all adaptors for multi-modal inputs use Sigmoid-Weighted Linear Unit (SiLU) activation.

**Training Configuration of Base Policy.** For each task, we collect 200 demonstrations with tele-operation to train a corresponding model from scratch. Similar to the implementation of RDT-1B, the vision encoder is frozen during training. The hyper-parameters for training are listed in Table 12.

**Implementation of Safety Alignment.** Following the implementation of aligning RDT-1B in simulation, we use LoRA to fine-tune the Transformer backbone with the same configuration as RDT-1B for simulation, and other hyperparameters are elaborated in Table 12 as well.

## E. Additional Results

### E.1. Qualitative Results for Simulation Evaluation

We further present qualitative results for three simulated tasks in Fig. 11. The base DP policy (top) exhibits typical safety violations, including gripper poking, behavior misalignment, and incorrect gripper rotation, which often lead to failures. After post-training with PACT (bottom), these failure modes are largely corrected, yielding safer interactions and higher task

Table 12. Hyper-parameter settings for real world evaluation.

Parameter	Value
Optimizer	AdamW
Optimizer betas	[0.9, 0.9999]
<b>Base policy training</b>	
batch size	$64 \times 8$
Learning rate	$1.0 \times 10^{-4}$
Learning-rate scheduler	Constant
Warmup steps	500
Training iterations	50,000
Training GPU	DGX B200
<b>Safety alignment distillation</b>	
batch size	$32 \times 8$
Number of rollouts	40
Learning rate	$1.0 \times 10^{-4}$
Learning-rate scheduler	Constant
Training epochs	8
Inner epochs	25
Training GPU	DGX B200

completion as demonstrated in Table 1.

## E.2. Effect of Base Policy Quality on Post-Training Gains

The magnitude of improvement depends on the competence of the initial policy. We observe smaller absolute gains for weaker base policies (e.g., DP3), whose failures are often dominated by missing core skills, leaving limited scope for safety-aware refinement. In contrast, stronger policies (e.g., RDT-1B) benefit more consistently, suggesting that PACT primarily serves as a post-training safety and reliability enhancer that is most effective when the base policy already possesses reasonable task understanding, and the remaining errors stem from unsafe dynamics or marginal failure cases. Improvements are also more pronounced on difficult, precision-critical tasks (e.g., Pour Water and Stack Blocks), where minor control errors can readily trigger safety violations and thus strongly affect both Succ. and Safe. Conversely, for relatively simple or near-saturated tasks (e.g., Handover Apple), baseline performance is already high and PACT yields only marginal gains. Notably, in tasks with extremely tight safety margins (e.g., stacking), Safe may remain a bottleneck even as Succ. improves, reflecting the strictness of the safety criteria rather than incomplete task execution.

## E.3. Results of collapsed On-policy baselines

We also attempted to apply representative on-policy baselines, AWR (Peng et al., 2019), RFT (Gilks & Wild, 1992), and QSM (Psenka et al., 2024) to our setting. However, all three methods exhibited pronounced training instability despite extra efforts for training stability introduced, including soft updates (Haarnoja et al., 2018) and a KL regularizer (Schulman et al., 2017). As shown in Fig. 12, they may reach moderate success rates in the initial iterations, yet performance quickly degrades with continued training and ultimately collapses to low (and in some cases near-zero) success. Consequently, we omit these on-policy baselines from the main comparisons since their collapse precludes a meaningful evaluation.

We further analyze the collapse of these on-policy baselines and find that it is consistent with a self-reinforcing distribution-drift failure mode. AWR and RFT are weighted imitation learning methods; RFT can be viewed as a binarized variant of AWR that retains samples above an advantage threshold. In our safety-alignment regime, advantages estimated from trajectory costs are noisy, so exponentiation (AWR) or hard selection (RFT) induces severe weight degeneracy and sharply reduces the effective sample size. As a result, updates concentrate on a small, non-representative subset of rollouts, accelerating support

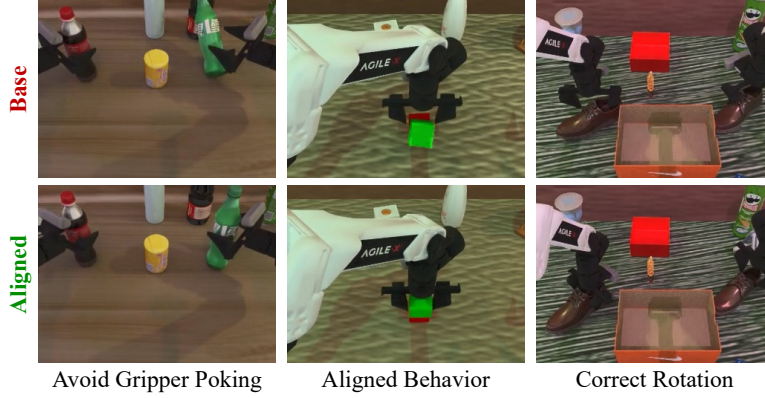


Figure 11. **Qualitative results in simulation.** Base policy (top) vs. Aligned Policy (bottom) on *Pick Dual Bottles*, *Stack Blocks*, and *Place Dual Shoes* (left-to-right). PACT reduces unsafe behaviors including poking, misalignment in both position and gripper pose.

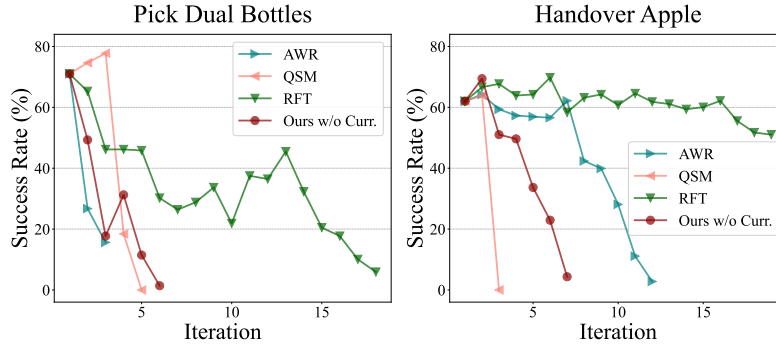


Figure 12. **Success rate (%) of collapsed on-policy baselines.** Success rate versus training iteration for three on-policy methods (AWR, QSM, and RFT) and our ablation *Ours w/o Curr.* on *Pick Dual Bottles* (left) and *Handover Apple* (right). Although some methods achieve moderate success in early iterations, performance is highly unstable: AWR and QSM rapidly collapse to near-zero success, and RFT shows substantial degradation over training (especially on *Pick Dual Bottles*), preventing a meaningful comparison in the main results. Meanwhile, *Ours w/o Curr.* also degrades and eventually gets near zero success.

shrinkage and triggering catastrophic forgetting of useful skills. Through closed-loop rollouts, this policy shift compounds over iterations and quickly drives the state distribution into safe-but-unproductive regions, where task-utility supervision is unavailable to recover competence, rendering the degradation effectively irreversible. By comparison,  $\text{PLD}_{\text{online}}$  mitigates collapse by introducing intervention of the implicit safe teacher at the end of each rollout, and provides external safe and useful behaviors as a supervision signal to preserve task-relevant behavior.

Meanwhile, QSM exhibits a related instability: its learned  $Q$ -function is queried on noisy diffusion-time inputs, yielding high-variance gradients and numerically brittle supervision. Moreover, QSM can be interpreted as a special case of direct distillation with a normalized teacher score. Consider the normalized distillation objective

$$\mathcal{J}(\theta) = \mathbb{E}_{t, \epsilon, (s, \mathbf{a}) \sim d^{\pi_{\theta}}} \left\| \epsilon_{\theta}(\mathbf{a}_t, \mathbf{s}, t) - \left( \epsilon_{\phi}(\mathbf{a}_t, \mathbf{s}, t) - \frac{\sum_{k=1}^m \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t) / \sigma_t}{\sum_{k=1}^m \lambda_k} \right) \right\|^2 \quad (30)$$

When the multipliers are uniformly large,  $\lambda_k \rightarrow \infty$  with a comparable scale, the base-score term vanishes after normalization and the target reduces to a (weighted) average constraint-gradient direction, giving

$$\begin{aligned} \mathcal{J}(\theta) &= \mathbb{E}_{t, \epsilon, (s, \mathbf{a}) \sim d^{\pi_{\theta}}} \left\| \epsilon_{\theta}(\mathbf{a}_t, \mathbf{s}, t) - \left( \epsilon_{\phi}(\mathbf{a}_t, \mathbf{s}, t) - \frac{\sum_{k=1}^m \lambda_k \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t) / \sigma_t}{\sum_{k=1}^m \lambda_k} \right) \right\|^2 \\ &= \mathbb{E}_{t, \epsilon, (s, \mathbf{a}) \sim d^{\pi_{\theta}}} \left\| \epsilon_{\theta}(\mathbf{a}_t, \mathbf{s}, t) + \frac{1}{m} \sum_{k=1}^m \nabla_{\mathbf{a}_t} c_{k,t}(\mathbf{s}, \mathbf{a}_t) / \sigma_t \right\|^2 \\ &= \mathcal{J}_{\text{QSM}}(\theta) \end{aligned} \quad (31)$$

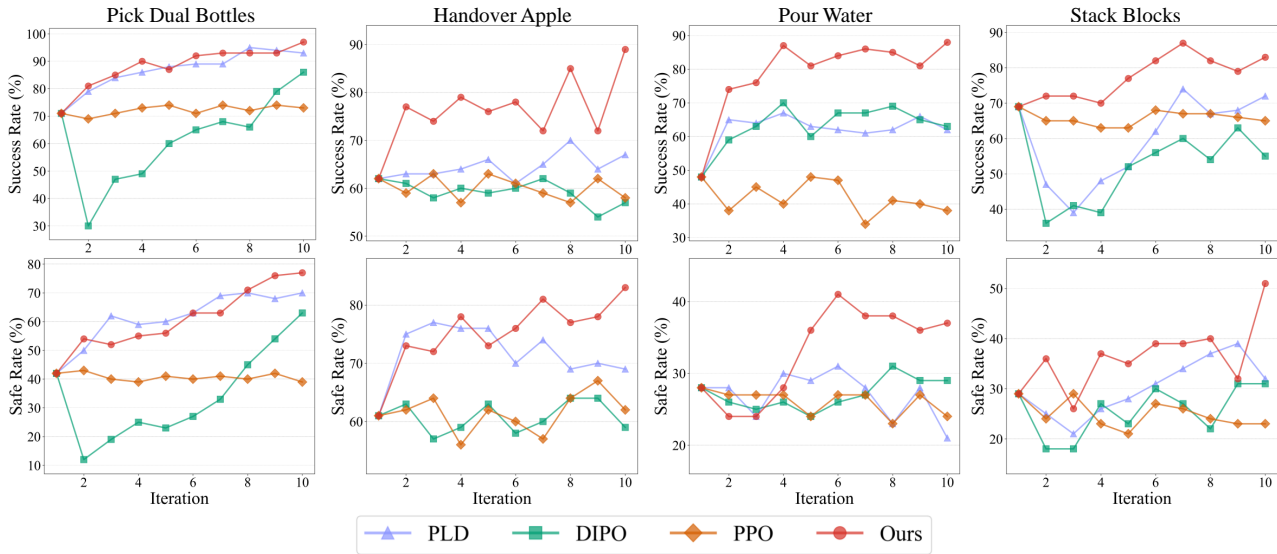


Figure 13. **Training efficiency comparison with on-policy baselines.** Success Rate (top) and Safe Rate (bottom) over training iterations on four tasks: *Pick Dual Bottles*, *Handover Apple*, *Pour Water*, and *Stack Blocks*. Compared with PLD, DIPO, and PPO, our method reaches higher performance in fewer iterations and yields more stable training curves, achieving the best final Success/Safe Rates across tasks.

Consequently, excessively large multipliers correspond to aggressive constraint enforcement and can induce collapse, consistent with our ablations in Sec. 4. In addition, mismatched numerical scales between predicted scores and cost gradients can further destabilize optimization, exacerbating QSM’s sensitivity.

#### E.4. Detailed training efficiency comparison with on-policy baselines

We compare the training dynamics of PACT with three on-policy baselines in Fig. 13 on four tasks: *Pick Dual Bottles*, *Handover Apple*, *Pour Water*, and *Stack Blocks*. Specifically, we track both task success rate and safety rate throughout training. PACT consistently exhibits superior training efficiency, achieving higher Success and Safe Rates earlier and attaining the best final performance across all four tasks.

#### E.5. More Ablation Studies

**Curriculum distillation.** We initially investigate the effect of curriculum safety constraints during training, which explicitly define a smooth transition between the base and the aligned policy by introducing a guidance schedule. We compare direct constraint distillation (Eq. (9)) against progressively tightening the constraint multipliers using a simple linear schedule  $\eta_i = i/N$  as described in Eq. (11). We find that direct distillation frequently induces abrupt distributional shifts, leading to the collapse of the action support and noticeable performance drops consistent with catastrophic forgetting (*Ours w/o Curr.* in Fig. 12). In contrast, curriculum-based enforcement yields smoother optimization dynamics and maintains continuous improvement.

**Further discussion on Lagrange multipliers  $\lambda$ .** The Lagrange multiplier  $\lambda$  controls the strength of constraint guidance during alignment. As shown in Fig. 6 with grid search, the guidance strength exhibits a clear trade-off between constraint enforcement and policy stability. When the multiplier is excessively large (e.g.,  $\times 5$  or  $\times 10$ ), the policy collapses across most tasks, failing to make meaningful progress, which corresponds to direct distillation that enforces constraints and leads to irreversible OOD collapse. In contrast, extremely small multipliers ( $\times 0.1$ ) lead to under-enforced constraints; the alignment process becomes inefficient and yields limited safety improvements despite being stable. The optimal performance consistently emerges near the turning point of this trade-off, which enables simple hyperparameter selection. These results validate the necessity of curriculum distillation and highlight the robustness of PACT within a reasonable and practical multiplier range.

**Further discussion on cost distillation with few diffusion steps.** As shown in Fig. 7, increasing the number of guided denoising steps does not monotonically improve performance. Our experiments demonstrate that injecting constraint guidance for only a few denoising steps at the late stage (e.g.,  $t_c = 0.03$ ) is sufficient to achieve strong alignment, whereas more aggressive distillation ( $t_c \geq 0.1$ ) often degrades both task success and safety performance.

This behavior can be attributed to the noise structure of diffusion models. Early denoising steps correspond to high-noise regimes, particularly when  $t_c \rightarrow 1$ , where sampled actions are dominated by near-Gaussian randomness. Applying safety costs at this stage yields gradients that are weakly correlated with the final action realization and may significantly distort the original diffusion sampling distribution, thereby impairing stable and meaningful policy optimization. These findings empirically justify our efficient design choice of late-stage, few-step constraint distillation, which aligns well with prior observations in guided diffusion while preserving the integrity of the pretrained policy distribution.