


STARE-VLA: Progressive Stage-Aware Reinforcement for Fine-Tuning Vision-Language-Action Models

Feng Xu^{*,1,3}, Guangyao Zhai^{*,1}, Xin Kong^{†,2}, Tingzhong Fu³, Daniel F.N. Gordon³, Xueli An³ and Benjamin Busam¹

¹Technical University of Munich, ²Imperial College London, ³Munich Research Center, Huawei Technologies

*Equal Contributions, [†]Corresponding Authors

✉ EAL.Feng.Xu@tum.de, x.kong21@imperial.ac.uk,  STARE-VLA

Abstract | Recent advances in Vision-Language-Action (VLA) models, powered by large language models and reinforcement learning-based fine-tuning, have shown remarkable progress in robotic manipulation. Existing methods often treat long-horizon actions as linguistic sequences and apply trajectory-level optimization methods such as Trajectory-wise Preference Optimization (TPO) or Proximal Policy Optimization (PPO), leading to coarse credit assignment and unstable training. However, unlike language, where a unified semantic meaning is preserved despite flexible sentence order, action trajectories progress through causally chained stages with different learning difficulties. This motivates progressive stage optimization. Thereby, we present **Stage-Aware Reinforcement (STARE)**, a module that decomposes a long-horizon action trajectory into semantically meaningful stages and provides dense, interpretable, and stage-aligned reinforcement signals. Integrating STARE into TPO and PPO, we yield Stage-Aware TPO (**StA-TPO**) and Stage-Aware PPO (**StA-PPO**) for offline stage-wise preference and online intra-stage interaction, respectively. Further building on supervised fine-tuning as initialization, we propose the Imitation → Preference → Interaction (**IPI**), a serial fine-tuning pipeline for improving action accuracy in VLA models. Experiments on SimplerEnv and ManiSkill3 demonstrate substantial gains, achieving state-of-the-art success rates of 98.0% on SimplerEnv and 96.4% on ManiSkill3 tasks.

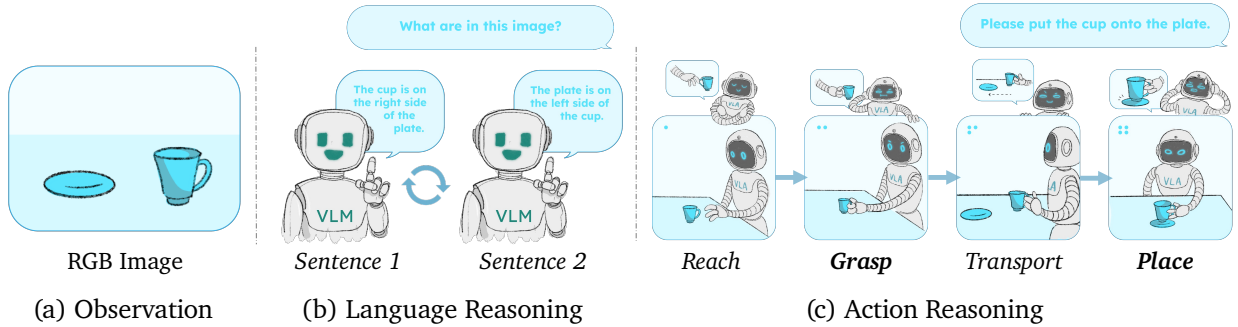


Figure 1 | Language Reasoning vs. Action Reasoning. Given an RGB image as the observation (a), the language model (b) is asked to describe the content in the image, and produces *Sentence 1* and *Sentence 2*. These sentences are flexibly ordered and jointly contribute to the global meaning required to answer the question. In contrast, the VLA model (c), when instructed to place the cup onto the plate, generates an action trajectory composed of semantically meaningful stages (*Reach*→*Grasp*→*Transport*→*Place*), which follow a strict order and vary in difficulty (with the more challenging stages shown in bold).

1. Introduction

Large-scale Vision-Language-Action (VLA) models [Zitkovich et al., 2023; Ghosh et al., 2024; Kim et al., 2024; Black et al., 2024, 2025] have recently emerged as powerful generalist policies for robotic manipulation. These models unify image, language, and action modalities within a single architecture, enabling robots to interpret multimodal inputs and generate executable actions. Pretrained on massive-scale multimodal data [O’Neill et al., 2024; Walke et al., 2023], VLA models provide strong priors that can be efficiently adapted to diverse downstream tasks through fine-tuning, avoiding the need for retraining from scratch.

Recent development of large-scale VLA models has been rapidly driven by the success of vision-language models (VLMs) and large language models (LLMs), as their output, i.e., action trajectories and sentences, can both be represented as sequential data [Zitkovich et al., 2023; Ghosh et al., 2024; O’Neill et al., 2024]. Consequently, many developed fine-tuning techniques, such as supervised fine-tuning (SFT), Reinforcement Learning from Human Feedback (RLHF) [Ouyang et al., 2022], direct preference optimization (DPO) [Rafailov et al., 2023], Proximal Policy Optimization (PPO) [Schulman et al., 2017], and Group Relative Policy Optimization (GRPO) [Guo et al., 2025a], have been straightforwardly adopted for VLA models. However, directly applying these methods to fine-tune on the whole action trajectories remains cumbersome and often inefficient, as the large optimization space makes credit assignment across long-horizon trajectories highly ambiguous. Unlike language reasoning, where optimization depends on a holistic understanding of sentences without strict ordering, an action trajectory naturally decomposes into semantically distinct stages that are causally chained and vary in difficulty. For example, as a pick-and-place task illustrated in Figure 1, *Reach* must precede *Grasp*, which in turn precedes *Transport* and *Place*. *Reach* and *Transport* are relatively easy with simple optimization objectives, while *Grasp* and *Place* are more challenging as they require precise geometric constraints. Overall task success hinges on correct progression through all stages. This fundamental characteristic motivates stage-aware objectives rather than monolithic trajectory-level optimization, which remains the predominant paradigm in current VLA fine-tuning.

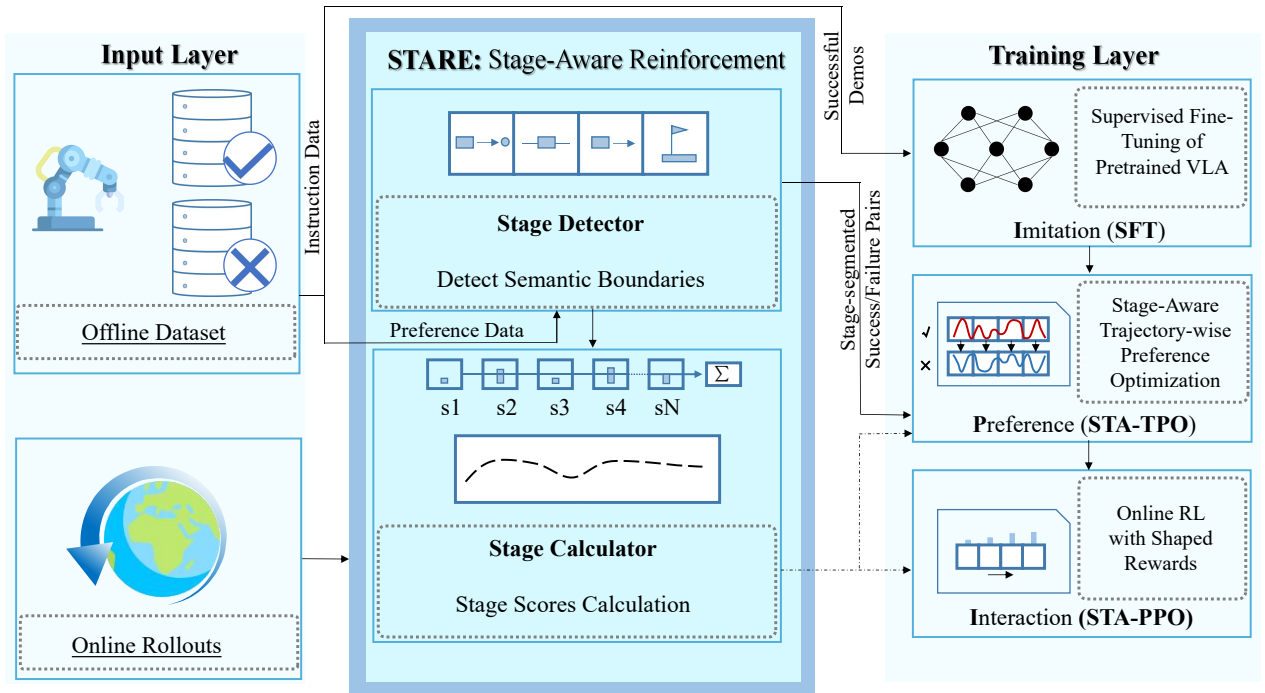


Figure 2 | Overview of the STARE Framework and Its Integration into the IPI Training Pipeline.

In this paper, we design *Stage-Aware Reinforcement* (STARE), a plug-in module that decomposes action trajectories into progressive stages with dense reward signals based on task-specific semantics. Given a trajectory, either in the collected data or during the model’s rollout, STARE employs a stage separator to identify *when* stage transitions occur, based on the translation and orientation of an end-effector. A stage calculator computes a stage cost and per-step rewards to evaluate *how well* each stage is executed. In this way, STARE not only annotates stage-wise actions but also assesses partial successes and failures within a trajectory. We leverage STARE for offline fine-tuning via *Stage-Aware Trajectory-Wise Preference Optimization* (STA-TPO), which constructs pairwise preferences at the stage level. By incorporating stage costs, STA-TPO propagates precise gradient signals to specific action stages, enabling progressive learning and credit assignment not only between success and failure but also among varying degrees of success. For online fine-tuning, we introduce STARE to *Stage-Aware Proximal Policy Optimization* (STA-PPO), which reshapes sparse terminal rewards to dense interaction rewards. By providing this progressive feedback, STA-PPO stabilizes intra-stage updates, especially for complex manipulation tasks that require dense guidance. Conceptually, STA-TPO and STA-PPO are reminiscent of curriculum learning [Bengio et al., 2009], where training is organized along an ordered sequence of subtasks to ease optimization and improve generalization. However, unlike conventional curricula that progress strictly from easy to hard, our stage-aware design enforces semantic continuity across stages, ensuring that optimization respects causal dependencies in stages.

To sufficiently fine-tune a pre-trained VLA model with STA-TPO and STA-PPO, we integrate these two algorithms with SFT as an initialization into a serial tri-step fine-tuning pipeline, *Imitation* \rightarrow *Preference* \rightarrow *Interaction* (IPI). The IPI framework first finetunes a VLA model from expert demonstrations via SFT, then further optimizes it according to offline stage-aware preferences using STA-TPO, and finally refines it based on stage-aware interaction in online environments using STA-PPO. In contrast to existing VLA fine-tuning strategies, IPI offers two key advantages: first, it explicitly models the multi-stage structure of robot trajectories, enabling more precise credit assignment rather than monolithic trajectory-level optimization [Zhang et al., 2024]. Second, compared to other methods that treat offline and online fine-tuning as disjoint processes [Zhang et al., 2024; Lu et al., 2025a; Chen et al., 2025a], IPI unifies them under a single framework, enabling stage-wise preference alignment and intra-stage interactions. Extensive experiments show that IPI not only improves in-distribution success rates but also significantly enhances out-of-distribution generalization, underscoring the importance of multi-stage reward design in VLA fine-tuning.

Our contributions are summarized as: **(i)** We design STARE, a rule-based module that decomposes trajectories into semantically meaningful stages, enabling fine-grained supervision beyond trajectory-level signals. **(ii)** Based on (i), we propose stage-aware fine-tuning methods: STA-TPO for offline stage-wise preference alignment and STA-PPO for online intra-stage interaction, both providing more precise credit assignment and improved sample efficiency. **(iii)** We unify supervised fine-tuning, STA-TPO, and STA-PPO into IPI, a serial tri-step pipeline for fine-tuning VLA models, and validate it on the benchmarking frameworks SimplerEnv and ManiSkill3, showing that IPI achieves state-of-the-art success rates.

2. Related Work

Long-horizon Robotic Manipulation Tasks in RL Long-horizon robotic manipulation involves completing a sequence of sub-tasks with frequent state and environment changes. Applying RL to such tasks is challenging due to sparse rewards, credit assignment, error accumulation, and high-dimensional state spaces. To address these issues, Plan-Seq-Learn [Dalal et al., 2024] leverages language models for high-level planning and RL for low-level control, enabling end-to-end execution from visual input to complex tasks. AC³ [Yang et al., 2025] learns continuous action chunks with intrinsic rewards to mitigate sparsity. DEMO³ [Escoriza et al., 2025] augments limited demonstrations with a world model and stage-wise dense rewards to improve sample efficiency. RoboHorizon [Chen et al., 2025b] employs LLMs to generate sub-goals and rewards, integrated with multi-view world models and planning to achieve high success rates. ARCH [Sun et al., 2025] combines high-level policy selection with a primitive skill library to tackle contact-rich assembly. SARM [Chen et al., 2025c] uses subtask-annotated reward modeling to filter demonstration quality, enabling robust long-horizon deformable object manipulation. Building on these inspirations that divide goals into sub-goals, we take a further step with stage-aware reinforcement, decomposing trajectories into semantically meaningful stages. This provides denser feedback and enables progressive optimization, making RL more effective for long-horizon VLA tasks.

RL Fine-tuning for LLMs RL fine-tuning is a widely used approach for aligning LLMs. The most prominent method is RLHF [Ouyang et al., 2022], where a reward model trained from human preference data guides algorithms such as policy gradient or PPO to align outputs with human expectations. While highly successful, RLHF suffers from costly data collection and unstable training. To mitigate these issues, DPO [Rafailov et al., 2023] eliminates the reward model by directly optimizing on preference comparisons, simplifying training and improving stability. Further variants such as RLAIIF [Lee et al., 2023] and RAFT [Dong et al., 2023] refine the framework. DeepSeek-R1 [Guo et al., 2025a] employs GRPO, which samples multiple responses per prompt and uses their relative performance within the group to compute advantages. As a subclass of LLMs, Large Reasoning Models (LRMs) [Zhang et al., 2025a] utilize Chain-of-Thought (CoT) [Wei et al., 2022] or Process Reward Model (PRM) [Lightman et al., 2023] for multi-step reasoning and face challenges akin to long-horizon VLA tasks, including sparse rewards and difficult credit assignment. This motivates our stage-aware reinforcement approach for VLA models.

RL Fine-tuning for VLAs Recent studies explore RL as a fine-tuning paradigm for VLA models. GRAPE [Zhang et al., 2024] adapts DPO [Rafailov et al., 2023] to trajectory-level preferences to propose TPO, while ConRFT [Chen et al., 2025a] alternates RL and SFT in real-world settings. ReinboT [Zhang et al., 2025b] designs dense rewards, and Guo et al. [2025b] propose an iterative SFT-RL pipeline to reduce instability and cost. RIPT-VLA [Tan et al., 2025] applies RLOO [Ahmadian et al., 2024] for online training, RL4VLA [Liu et al., 2025] studies RL-driven generalization, VLA-RL [Lu et al., 2025a] applies PPO, and RFTF [Shu et al., 2025] introduces value models for dense reward estimation. SimpleVLA-RL [Li et al., 2025a] extends verL to VLA models with GRPO-based online RL, demonstrating significant improvements in data efficiency, long-horizon task performance, and generalization across spatial, object, and task distributions. RLinf [Zang et al., 2025] provides a scalable and unified pipeline for VLA RL—combining rendering, inference, and training—to boost efficiency and performance. π_{RL} [Chen et al., 2025d] applied online RL fine-tuning for flow-based VLAs. $\pi_{0.6}^*$ [Intelligence et al., 2025] uses RECAP [Intelligence et al., 2025] to fine-tune VLAs with advantage-conditioned policies, learning from autonomous experience and expert corrections. GR-RL [Li et al., 2025b] combines offline data filtering with distributional critics and online latent-

space RL for high-precision dexterous manipulation. Despite their promise, these methods typically optimize at the trajectory level, suffering from sparse rewards, coarse credit assignment, and difficult exploration in long-horizon manipulation. In contrast, our stage-aware RL decomposes trajectories into semantically meaningful stages and assigns stage-level rewards, providing denser, interpretable feedback and enabling progressive optimization for complex robotic tasks.

3. Preliminary

3.1. Problem Formulation

We consider a language-conditioned POMDP problem defined by the tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{L}, \mathcal{R}, \gamma\}$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is the dynamic function, \mathcal{L} is the space of language instruction, $\mathcal{R} : \mathcal{S} \times \mathcal{L} \rightarrow \mathbb{R}$ is the reward function, and γ is a scale factor with $0 < \gamma < 1$. The goal of a VLA model is to find a policy $\pi_\theta : \mathcal{S} \times \mathcal{L} \rightarrow \mathcal{A}$, which generates action trajectories maximizing the expected accumulated reward, or return for each task l , i.e. $\mathcal{R}(\pi, l) = \mathbb{E}_{a \sim \pi} [\sum_t \gamma^t r_t]$.

Fine-tuning a VLA model adapts a pre-trained π_θ to new tasks so that the resulting policy $\pi_{\theta'}$ maximizes expected return under the POMDP. This can be done through imitation for aligning with expert demonstrations, preference for refining trajectories via learned comparisons, or reinforcement optimizing long-term rewards.

3.2. Trajectory-Wise Preference Optimization (TPO)

Direct Preference Optimization (DPO) [Rafailov et al., 2023] is a recent fine-tuning technique for large language models that directly aligns a policy with preference data, bypassing explicit reward modeling. Extending this idea to fine-tuning VLA models yields TPO [Zhang et al., 2024]: the outputs are action trajectories $\tau = \{(s_t, a_t)\}_{t=1}^T$ rather than text sequences. TPO treats each trajectory as a single sequence and learns from pairwise comparisons of successful and failed trajectories (τ^+, τ^-) generated under the same instruction. The policy is updated to prefer τ^+ over τ^- by minimizing

$$L_{\text{TPO}}(\theta) = -\mathbb{E}_{(\tau^+, \tau^-)} \left[\log \sigma(\beta(q(\tau^+) - q(\tau^-))) \right], \quad (1a)$$

$$q(\tau) = \frac{1}{T} \sum_{t=1}^T \left(\log \pi_{\theta'}(a_t | s_t) - \log \pi_\theta(a_t | s_t) \right). \quad (1b)$$

where $\sigma(\cdot)$ is the sigmoid function, β controls the strength of preference alignment, $s_t \in \mathcal{S}$ and $a_t \in \mathcal{A}$ denote the environment state and action at timestep t , and $q(\cdot)$ measures the normalized log-likelihood ratio of a trajectory under policy $\pi_{\theta'}$ relative to π_θ . L_{TPO} is minimized when the model increases $q(\tau^+)$ relative to $q(\tau^-)$, i.e., when the likelihood of successful trajectories exceeds failed ones.

While TPO provides a direct mechanism to apply preference learning to long-horizon control, it suffers from credit assignment ambiguity: preferences are assigned to full trajectories, making it difficult to determine which specific stage contributed to the preference signal. Moreover, such a binary preference limits optimization to coarse distinctions between successful and failed rollouts, without capturing relative quality among partially successful trajectories. These limitations motivate (STA-TPO), which decomposes trajectories into stages and aligns hierarchical preferences at the stage level, enabling finer-grained optimization.

3.3. Proximal Policy Optimization (PPO)

PPO [Schulman et al., 2017] is one of the most widely used online reinforcement learning algorithms, known for its balance of sample efficiency and training stability. PPO improves policy gradient methods by introducing a clipped surrogate objective that prevents excessively large policy updates, thereby stabilizing training. Given an old policy π_θ , the clipped objective is

$$L_{\text{PPO}}(\theta) = \mathbb{E}_t \left[\min \left(p_t(\theta) \text{GAE}(r_t), \text{clip}(p_t(\theta), 1 - \epsilon, 1 + \epsilon) \text{GAE}(r_t) \right) \right], \quad (2)$$

where $p_t(\theta) = \pi_{\theta'}(a_t|s_t)/\pi_\theta(a_t|s_t)$ is the likelihood ratio between the new and old policies, and ϵ is a clipping parameter. $\text{GAE}(\cdot)$ is generalized advantage estimator [Schulman et al., 2015] that estimate the advantage value based on rewards r_t .

In the context of fine-tuning VLA models, PPO is commonly used to fine-tune policies with sparse r_t , but such signals often limit sample efficiency and provide insufficient guidance for complex, long-horizon tasks. This motivates STA-PPO, which integrates stage-aware reward shaping to transform sparse terminal rewards into dense progressive signals for more efficient fine-tuning.

4. Method

We begin by introducing *Stage-Aware Reinforcement* (STARE), which decomposes long-horizon action trajectories into semantically meaningful stages, each equipped with stage-wise costs and intra-stage rewards. Building on this foundation, we develop offline and online learning algorithms for progressive stage-aware fine-tuning: *Stage-Aware Trajectory Preference Optimization* (STA-TPO) and *Stage-Aware Proximal Policy Optimization* (STA-PPO). Finally, we integrate STA-TPO and STA-PPO with supervised fine-tuning (SFT) into a serial pipeline, *Imitation* \rightarrow *Preference* \rightarrow *Interaction* (IPI), to achieve sufficient fine-tuning of VLA models.

4.1. Stage-Aware Reinforcement (STARE)

We propose STARE, a module that decomposes long-horizon action trajectories into semantically meaningful stages defined by task-specific rules. STARE consists of two components: (i) a *stage separator*, which determines *when* stage transitions occur by detecting task-relevant events, and (ii) a *stage calculator*, which evaluates *how well* each stage is executed using stage-wise costs and dense intra-stage rewards.

Stage Separator Stage boundaries are determined by semantically meaningful manipulation events rather than arbitrary temporal cuts. Given the whole action trajectory τ , we intend to divide it into K stages by defining semantic boundaries and assigning each global timestep t a stage label $k \in \{1, \dots, K\}$. Following an event-driven rule, the entry condition of stage k coincides with the terminal condition of stage $k - 1$, ensuring progressive continuity across stages. For instance, a pick-and-place task can be separated into four stages: *Reach* \rightarrow *Grasp* \rightarrow *Transport* \rightarrow *Place*.

Stage segmentation thus reduces to detecting the onset of each stage based on geometric constraints, defined by thresholds δ_k on the translation and orientation signals of the end-effector. These thresholds set binary environment flags (e.g., grasped, on-target). For example: a *Reach* \rightarrow *Grasp* transition occurs when the end-effector contacts the object; *Grasp* \rightarrow *Transport* occurs when the grasped object is lifted above a small height threshold; *Transport* \rightarrow *Place* occurs when the object is within a distance margin of the goal position; and *Place* \rightarrow *Success* occurs when the object is released and remains stably in the goal region (see Supplementary Material C for other segmentation examples). Thereby,

we group steps with the same stage label into the k -th trajectory segment $\tau^{(k)} = \{(s_t, a_t) \mid g(t) = k\}_{t=1}^{T_k}$, where $s_t \in \mathcal{S}$, $a_t \in \mathcal{A}$, and T_k is the number of timesteps assigned to stage k . Here, $g : \mathbb{N} \rightarrow \{1, \dots, K\}$ is a stage assignment function mapping each timestep t to its corresponding stage index k . The full trajectory can then be expressed as the stage-wise decomposition $\tau \mapsto \{\tau^{(i)}\}_{i=1}^K$.

Stage Calculator Given the stage segments produced by the stage separator, the stage calculator computes both stage-wise costs and intra-stage dense rewards by measuring the relation between the end-effector and relevant targets in the environment. The specific forms of cost and reward depend on the goal of each stage. We illustrate with *Reach* as the k -th stage:

(i) *Stage cost aggregation.* We define the cost function $\ell_k(\cdot)$ as the mean Euclidean distance over T_k between the end-effector and the target object from start to the end of *Reach*:

$$\ell_k(\tau^{(k)}) = \frac{1}{T_k} \sum_{t=1}^{T_k} \|x_{ee}(t) - x_{obj}(t)\|_2, \quad (3)$$

where $x_{ee}(t) \in \mathbb{R}^3$ denotes the Cartesian position of the end-effector at time step t , and $x_{obj}(t) \in \mathbb{R}^3$ is the target position of the object. By definition, ℓ_k is a non-negative value measuring the deviation from the target: the better the $\tau^{(k)}$, the smaller the ℓ_k . Detailed cost functions for other stage categories are provided in the Supplementary Material C.

(ii) *Intra-stage reward shaping.* To provide dense guidance, we adopt potential-based reward shaping [Kim et al., 2025; Ng et al., 1999]. For active stage k , we define a per-timestep potential Φ_{k_t} that captures the normalized progress of state s_t . Specifically, for *Reach*, we use:

$$\Phi_{k_t}(s_t) = \sigma\left(1 - \frac{\|x_{ee}(t) - x_{obj}(t)\|}{d_k}\right), \quad (4)$$

where $\Phi_{k_t}(s_t) \in [0, 1]$, d_k is a normalization length scale, and $\sigma(\cdot)$ is a sigmoid function. This provides smooth shaping rewards that encourage the end-effector to progressively reach the target (see detailed potential functions for other stages in the Supplementary Material C). Based on Φ_{k_t} , the shaped reward r'_t augments the sparse reward r_t as:

$$r'_t = r_t + \gamma \Phi_{k_{t+1}}(s_{t+1}) - \Phi_{k_t}(s_t). \quad (5)$$

4.2. From STARE to STA-TPO

Unlike standard TPO [Zhang et al., 2024], which aggregates preferences only at the level of entire trajectories, STA-TPO leverages STARE to segment trajectories into progressive stages and perform stage-wise preference alignment. A detailed algorithm is shown in Algorithm 1. A pair comparison of stage samples $(\tau^{(k)+}, \tau^{(k)-})$ exists only when the previous stage $\tau^{(k-1)}$ has been successfully completed, ensuring progressive consistency across stages. In addition, the stage cost $\ell_k(\tau)$ is incorporated as a penalty term in (1b), transforming q into \hat{q} :

$$\hat{q}(\tau^{(k)}) = q(\tau^{(k)}) - \lambda \ell_k(\tau^{(k)}), \quad (6)$$

where λ is the penalty weight. The original objective \mathcal{L}_{TPO} in (1a) thereby extends to $\mathcal{L}_{\text{STA-TPO}}$. Compared to (1), which optimizes the model only with binary trajectory-level preferences (success vs. failure), \hat{q} introduces a hierarchical signal to $\mathcal{L}_{\text{STA-TPO}}$. Even among successful stages $\tau^{(k)+}$ across different trajectories, those with lower penalties $\ell_k(\tau^{(k)})$ yield higher \hat{q} , while less optimal stages receive lower \hat{q} . This design enables credit assignment not only between success and failure but also among varying degrees of success, thereby providing finer-grained supervision for learning optimal behaviors.

Algorithm 1 STA-TPO (offline)

Require: Preference pairs $\{(\tau^+, \tau^-)\}$ under same instruction; reference policy π_θ ; STARE; stage penalty weight λ ; temperature β ; learning rate η .

Ensure: Updated policy $\pi_{\theta'}$.

- 1: **while** not converged **do**
- 2: Sample minibatch $\mathcal{B} = \{(\tau^+, \tau^-)\}$.
- 3: **for all** $\tau \in \{\tau^+, \tau^-\}$ in \mathcal{B} **do**
- 4: **Stage segmentation & costs:** $\{\tau^{(k)}, \ell_k(\tau^{(k)})\}_{k=1}^K \leftarrow \text{STARE}(\tau)$.
- 5: **Stage scores:** For each k , first compute

$$q(\tau^{(k)}) \leftarrow \frac{1}{T_k} \sum_{t \in T_k} \left(\log \pi_{\theta'}(a_t | s_t) - \log \pi_\theta(a_t | s_t) \right). \quad (\text{cf. (1b)})$$

- 6: Then add stage costs: $\hat{q}(\tau^{(k)}) \leftarrow q(\tau^{(k)}) - \lambda \ell_k(\tau^{(k)})$.
- 7: **end for**
- 8: **Preference loss:** Compute

$$L_{\text{STA-TPO}} \leftarrow -\frac{1}{|\mathcal{B}|} \sum_{(\tau^+, \tau^-) \in \mathcal{B}} \frac{1}{K} \sum_{k=1}^K \log \sigma \left(\beta (\hat{q}(\tau^{(k)+}) - \hat{q}(\tau^{(k)-})) \right). \quad (\text{cf. (1a)})$$

- 9: **Policy update:** $\theta' \leftarrow \theta' - \eta \nabla_{\theta'} L_{\text{STA-TPO}}$.
 - 10: **end while**
-

4.3. From STARE to STA-PPO

For online RL fine-tuning, we integrate STARE directly into rollouts. The stage separator determines the stage transition online. At each time step within the stage, the stage calculator produces shaped reward r'_t , turning L_{PPO} in (2) into $L_{\text{STA-PPO}}$. Finally, policy parameters θ are updated by minimizing $L_{\text{STA-PPO}}$. By replacing r_t with r'_t , STA-PPO provides denser, stage-aligned feedback that accelerates policy learning in long-horizon, sparse-reward tasks. A detailed algorithm is shown in Algorithm 2.

4.4. STA-TPO And STA-PPO for Serial Fine-tuning

Existing works often apply offline preference-based optimization [Zhang et al., 2024] and online RL fine-tuning [Liu et al., 2025; Li et al., 2025a] separately. Besides, while we are now able to jointly address offline preference alignment and online reinforcement learning by STA-TPO and STA-PPO, a complete fine-tuning framework for VLA models must also incorporate imitation learning to initialize a strong policy prior.

Thereby, we propose Imitation \rightarrow Preference \rightarrow Interaction (IPI), a three-step fine-tuning pipeline. We first warm up the policy safely from demonstrations by SFT. Then we apply STA-TPO to offline refine the policy. Finally, we apply STA-PPO to further enhance robustness through online exploration. Thereby, IPI integrates supervised, preference-based, and exploration signals into a coherent progression, yielding more sample-efficient and more robust fine-tuning of VLA models.

5. Experiments

Benchmarks & Baselines. We evaluate our approach on two families of robotic manipulation environments, as shown in Figure 3. The first is **SimplerEnv** [Li et al., 2024a] with the WidowX

Algorithm 2 STA-PPO (online)

Require: Simulation Env; Behavior policy π_θ ; STARE; horizon T ; PPO epochs E ; discount γ ; GAE parameter; clip ϵ ; step size η .

Ensure: Updated policy $\pi_{\theta'}$.

```

1: while not converged do
2:   Rollout: collect  $\{(s_t, a_t, r_t, \log \pi_\theta(a_t|s_t))\}_{t=0}^{T-1}$  in Env.
3:   Online stage labels & potentials:
4:   for  $t = 0$  to  $T - 1$  do
5:     Detect current stage  $k = g(t)$  via stage separator in STARE (event rules).
6:     Compute potential  $\Phi_{k_t}(s_t)$  by the stage calculator in STARE.
7:   end for
8:   Shaped rewards:
9:   for  $t = 0$  to  $T - 1$  do
10:     $r'_t \leftarrow r_t + \gamma \Phi_{k_{t+1}}(s_{t+1}) - \Phi_{k_t}(s_t)$  ▷ Potential-based shaping
11:  end for
12:  for  $e = 1$  to  $E$  do ▷ PPO updates
13:    Compute ratio  $p_t(\theta') = \exp(\log \pi_{\theta'}(a_t|s_t) - \log \pi_\theta(a_t|s_t))$ .
14:    Interaction loss:

$$L_{\text{STA-PPO}} \leftarrow \mathbb{E}_t [\min(p_t(\theta') \text{GAE}(r'_t), \text{clip}(p_t(\theta'), 1 - \epsilon, 1 + \epsilon) \text{GAE}(r'_t))] . \quad (\text{cf. (2)})$$

15:    Update:  $\theta' \leftarrow \theta' + \eta \nabla_{\theta'} L_{\text{STA-PPO}}$ .
16:  end for
17: end while

```

arm, where we focus on the four canonical single-object tasks in the **SimplerEnv-WidowX** split. The second is **ManiSkill3** [Tao et al., 2025] with the Franka robot [Haddadin, 2024], including *StackCube* and three contact-rich tasks (*PushCube*, *PullCube*, and *LiftPegUpright*) to validate generality beyond pick-and-place and assess performance under challenging non-trivial manipulation. We compare against widely used VLA baselines (RT-1-X, Octo-Base/Small, RoboVLM, SpatialVLA), the strong offline preference fine-tuning method GRAPE [Zhang et al., 2024], and the RL fine-tuning baselines RL4VLA [Liu et al., 2025] and π_{RL} [Chen et al., 2025d]. For fairness, all methods fine-tune the OpenVLA-7B [Kim et al., 2024] and pi0.5_base [Black et al., 2025] backbone, and we additionally evaluate our proposed STA-TPO, STA-PPO, and the full **IPI**. We report average success over 300 evaluation episodes per method and setting. Unless otherwise stated, hyperparameters are shared across methods when applicable (detailed in Appendix A).

Main Results. We begin by presenting overall comparisons on two widely used families of manipulation benchmarks. Table 1 reports grasp and final success rates on the four representative **SimplerEnv-WidowX** tasks, while Table 2 reports results on selected **ManiSkill3-Franka** tasks including both stacking and contact-rich manipulation.

Across all benchmarks, existing VLA baselines exhibit limited performance (e.g., average success rates $< 60\%$). Recent RL fine-tuning approaches, such as RL4VLA [Liu et al., 2025] achieve strong results (92.6% on **SimplerEnv-WidowX**, 70.5% on **ManiSkill3**). Our proposed **IPI** further improves to 98.0% and 96.4%, outperforming prior state-of-the-art methods by +5.4 and +25.9 points, nearly solving these benchmark tasks. Flow-based VLA models such as $\pi_{0.5}$ also benefit substantially from our stage-aware reinforcement design: integrating STARE consistently boosts $\pi_{0.5}$ ’s performance across all tasks, outperforming its original flow-matching baseline by a large margin. Our **IPI** is a fully

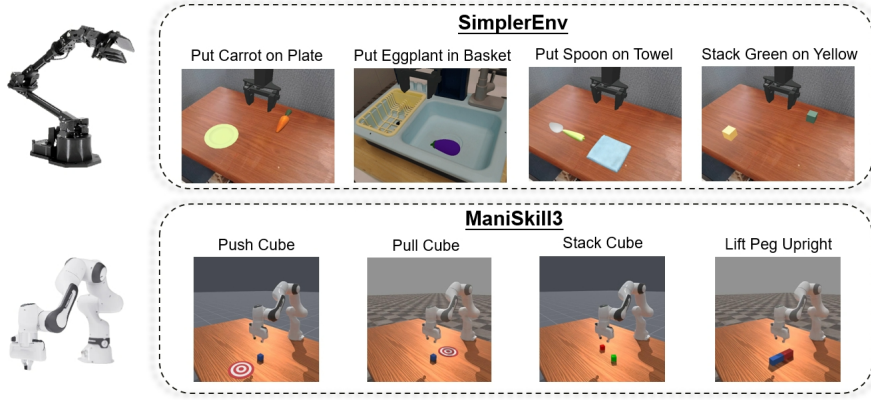


Figure 3 | Two simulated benchmarks. We show experiment setups and example tasks involved.

implemented and executed pipeline obtained from actual end-to-end runs, demonstrating that each stage can be integrated seamlessly and that the complete framework achieves the strongest overall performance.

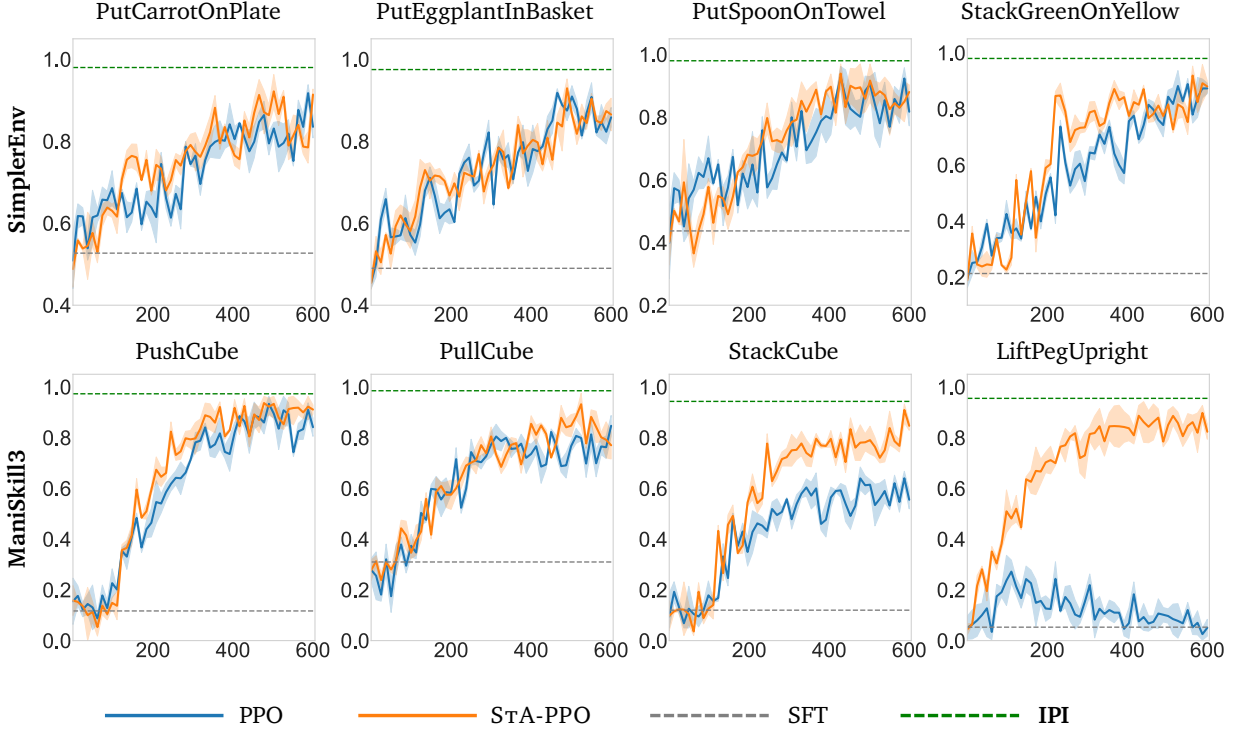


Figure 4 | Comparison of learning curves across eight representative tasks from SimplerEnv and ManiSkill3. The y-axis denotes the success rate, and the x-axis shows the interaction environment steps (in thousands). We compare SFT, PPO, ST-A-PPO, and our full IPI method. While PPO improves over SFT, it often stagnates in high-precision (e.g., *StackCube*) or contact-rich settings (e.g., *LiftPegUpright*). ST-A-PPO accelerates convergence and achieves higher asymptotic performance by leveraging stage-aware signals. Notably, the most challenging tasks, *LiftPegUpright* and *StackCube*, show the clearest benefit, highlighting the importance of incorporating stage-awareness for solving complex tasks.

PPO vs ST-A-PPO While PPO improves over SFT, it often stagnates on tasks requiring high precision or contact-rich interactions. In contrast, ST-A-PPO consistently accelerates convergence and achieves higher asymptotic performance by leveraging stage-aware signals. Figure 4 presents results across eight representative tasks from SimplerEnv-WidowX and ManiSkill3. The most challenging

Table 1 | Evaluation on SimplerEnv with WidowX Robot tasks. We report the final success rate and grasp success rate shown in parentheses (Success % (Grasp %)). Our method, IPI, uses RL fine-tuning after an initial SFT phase. The ‘(+X%)’ indicates the improvement over a relevant baseline.

Methods	Robotic Task				Avg. Success Rate (%)
	Put Spoon on Towel	Put Carrot on Plate	Stack Green on Yellow	Put Eggplant in Basket	
<i>Other Methods</i>					
RT-1-X [Brohan et al., 2023]	0.0 (16.7)	4.2 (20.8)	0.0 (8.3)	0.0 (0.0)	1.1
Octo-Base [Ghosh et al., 2024]	12.5 (34.7)	8.3 (52.8)	0.0 (31.9)	43.1 (66.7)	16.0
Octo-Small [Ghosh et al., 2024]	47.2 (77.8)	9.7 (27.8)	4.2 (40.3)	56.9 (87.5)	30.0
RoboVLM [Li et al., 2024b]	20.8 (37.5)	25.0 (33.3)	8.3 (8.3)	0.0 (0.0)	13.5
SpatialVLA [Qu et al., 2025]	20.8 (25.0)	20.8 (41.7)	25.0 (58.3)	70.8 (79.2)	34.4
SOFAR [Qi et al., 2025]	58.3 (62.5)	66.7 (75.0)	70.8 (91.7)	37.5 (66.7)	58.3
<i>OpenVLA-7B Based Methods</i>					
SFT	43.7(70.3)	52.7(74.7)	21.3(59.0)	49.0(67.3)	41.7
GRAPE [Zhang et al., 2024]	44.3(72.0)	55.0(85.3)	22.7(53.3)	53.7(78.7)	43.9
SFT → S _T A-TPO(STARE)	51.0(85.7)	57.3(82.3)	43.7(78.3)	54.3(85.7)	51.6 (+17.7)
RL4VLA [Liu et al., 2025]	93.0 (98.3)	91.3 (96.7)	92.0 (97.0)	93.7 (98.3)	92.5
SFT → S _T A-PPO(STARE)	94.3 (97.7)	95.3 (99.0)	93.7 (98.3)	95.0 (98.7)	94.6 (+2.1)
IPI (STARE)	98.0 (99.0)	98.5 (99.5)	98.0 (99.0)	97.5 (99.0)	98.0
<i>Pi0.5 Based Methods</i>					
SFT	49.3(85.3)	64.7(89.3)	44.7(76.0)	69.7(92.3)	57.1
GRAPE [Zhang et al., 2024]	48.0(78.7)	59.3(88.3)	48.3(69.7)	58.7(80.3)	53.6
SFT → S _T A-TPO(STARE)	54.0(83.7)	65.3(80.3)	54.0(70.7)	68.7(83.7)	60.5 (+6.9)
π_{RL} [Chen et al., 2025d]	82.7 (98.3)	97.3 (99.3)	83.3 (97.3)	55.0 (69.7)	79.6
SFT → S _T A-PPO(STARE)	90.7 (98.7)	97.7 (99.0)	85.7 (97.3)	63.7 (85.7)	84.5 (+4.9)
IPI (STARE)	95.7 (99.3)	98.7 (99.7)	93.0 (98.0)	78.7 (91.0)	91.5

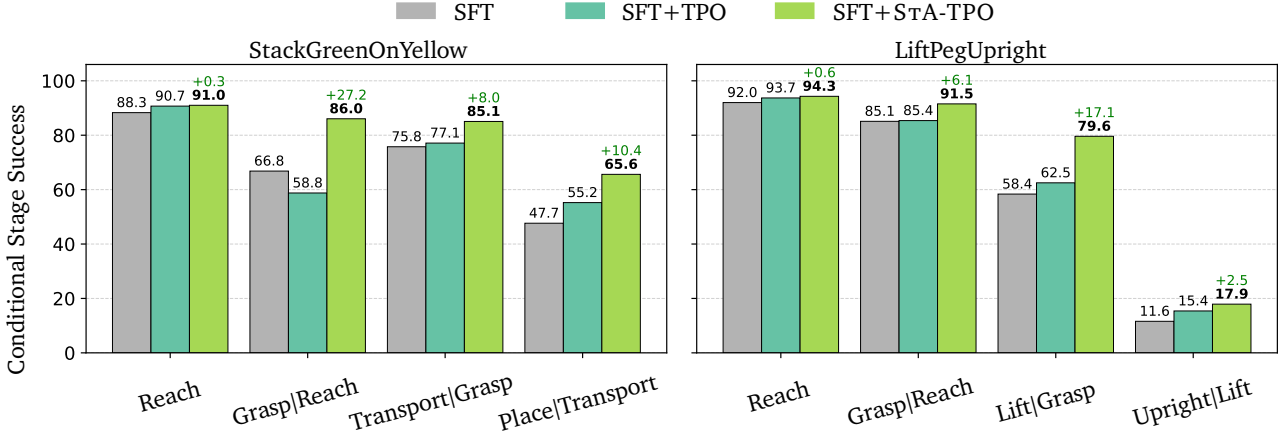
tasks—*LiftPegUpright* and *StackGreenOnYellow*—exhibit the largest performance gaps, underscoring the importance of incorporating stage-aware signals in long-horizon, precision-critical manipulation. By comparison, for short-horizon pick-and-place tasks (e.g., *PutCarrotOnPlate*, *PutEggplantInBasket*) or simple push-and-pull tasks, PPO and STA-PPO achieve similar final success rates, with STA-PPO mainly contributing faster convergence and reduced variance. Overall, these results suggest that stage-aware guidance is particularly crucial when strict alignment accuracy or multi-stage coordination is required, whereas simpler tasks can often be solved effectively with standard reinforcement learning.

After the benchmark-level comparison in Tables 1 and 2, we note that while the overall improvements of STA-PPO and STA-TPO over prior baselines are consistent, the performance gap is most pronounced on two tasks: (1) Cube stacking tasks from both the environments, which requires precise alignment in placing stage, and (2) *LiftPegUpright* from ManiSkill3, which demands accurate orientation control after lifting. To better understand where these gains originate, we decompose trajectories into semantic stages and evaluate **conditional stage success** ($P(\text{stage}_k | \text{stage}_{k-1})$), which measures how reliably a policy completes a stage given that all previous stages have been successful.

Figure 5 shows STA-TPO provides clear advantages over TPO, with the largest improvements appearing in the **grasp**, and **place** and **upright** stages. These stages are particularly decisive for the final outcome, explaining why the overall success rate improvements are disproportionately large for these two tasks.

Table 2 | Evaluation on selected ManiSkill3 Franka tasks. OpenVLA-7B and Pi0.5 based methods use 100 trajectory samples for SFT and 50 trajectory preference pairs for TPO and StA-TPO, detailed in Appendix B.4.

Methods	Robotic Task				Avg. Success Rate (%)
	Stack Cube	Push Cube	Pull Cube	LiftPeg Upright	
<i>Other Methods</i>					
Octo (fine-tuning) [Ghosh et al., 2024]	1.0	67.0	90.0	0.0	39.5
SmolVLA (fine-tuning) [Shukor et al., 2025]	12.7	86.3	90.7	16.3	51.5
RoboFAC-7B [Lu et al., 2025b]	85.5	80.4	80.7	84.0	82.7
<i>OpenVLA-7B Based Methods</i>					
SFT	12.0	11.7	31.0	5.3	15.0
GRAPE [Zhang et al., 2024]	15.7	13.3	35.3	7.7	18.0
SFT→StA-TPO (STARE)	19.3	16.0	35.7	12.3	20.8 (+2.8)
RL4VLA [Liu et al., 2025]	64.0	95.7	90.3	32.0	70.5
SFT→StA-PPO (STARE)	92.7	96.0	95.3	89.7	93.4 (+22.9)
IPI (STARE)	94.3	97.3	98.5	95.5	96.4
<i>Pi0.5 Based Methods</i>					
SFT	26.3	18.3	43.0	10.7	25.6
GRAPE [Zhang et al., 2024]	22.7	16.3	45.0	6.3	22.6
SFT→StA-TPO (STARE)	28.0	22.3	44.7	15.7	27.7 (+5.1)
π_{RL} [Chen et al., 2025d]	72.3	96.7	93.3	58.0	80.1
SFT→StA-PPO (STARE)	80.7	98.0	92.7	75.3	86.7 (+6.6)
IPI (STARE)	84.3	99.3	95.0	80.7	89.9

**Figure 5 | Offline Stage-wise ablation on two tasks.** We report stage completion rates (%) for *StackGreenonYellow* (SimplerEnv) and *LiftPegUpright* (ManiSkill3). Compared with TPO, StA-TPO achieves significant gains, particularly in the **grasp** and **place/upright** stages, which are critical for final success.

Ablations study To further dissect the contributions of stage-aware reinforcement, we conduct a stage toggle ablation where the STARE signal is selectively removed at different phases of the manipulation in StA-PPO. As shown in Figure 6, disabling STARE at early stages (e.g., reach or grasp) only leads to moderate drops, since later corrective actions can partially recover performance. In contrast, removing STARE at the final precision-critical phases (e.g., **Place** in stacking and **Upright** in peg lifting) causes the largest degradation, reducing success rates by more than 20%. This analysis highlights that STARE guidance is especially valuable at stages where geometric accuracy and stability

directly determine task completion.

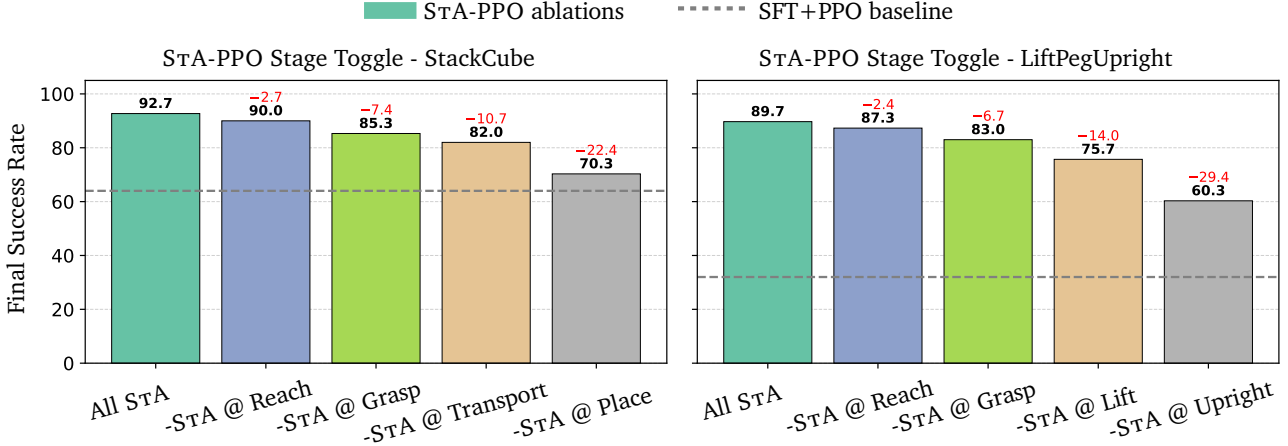


Figure 6 | Stage toggle ablation of STA-PPO. We evaluate the effect of selectively disabling stage-aware reinforcement signals on two representative tasks: *StackCube* (ManiSkill3) and *LiftPegUpright* (ManiSkill3). The **All STA (STARE)** setting achieves the best performance, while disabling critical stages (**Place** in stacking, **Upright** in peg lifting) causes the largest performance drops.

6. Conclusion

We presented *Stage-Aware Reinforcement* (STARE), a plug-in module that decomposes trajectories into semantically meaningful stages and provides stage-level reinforcement signals. Building on this, we introduced Stage-Aware TPO (STA-TPO) and PPO (STA-PPO) for offline stage-wise preference alignment and online intra-stage interaction, and integrated them with supervised fine-tuning into the *Imitation* \rightarrow *Preference* \rightarrow *Interaction* (IPI) pipeline. Experiments on SimplerEnv and ManiSkill3 demonstrate that IPI achieves state-of-the-art success rates.

Extensive experiments on both OpenVLA and Pi0.5 backbones, across SimplerEnv and ManiSkill3, show that IPI consistently yields substantial gains and achieves state-of-the-art success rates. These results demonstrate that stage-aware credit assignment is a key missing ingredient in current VLA reinforcement learning, enabling more stable optimization, faster convergence, and markedly stronger long-horizon performance. We believe this stage-centric perspective opens up promising directions for scalable, interpretable, and robust robotic learning.

References

- Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, 2023.
- Dibya Ghosh, Homer Rich Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, et al. Octo: An open-source generalist robot policy. In *Robotics: Science and Systems*, 2024.
- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, et al. Openvla: An open-source vision-language-action model. In *Conference on Robot Learning*, 2024.
- Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. $\pi 0$: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, et al. $\pi 0.5$: a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- Homer Rich Walke, Kevin Black, Tony Z. Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, Abraham Lee, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning*, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 2022.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 2023.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025a.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *International Conference on Machine Learning*, 2009.
- Zijian Zhang, Kaiyuan Zheng, Zhaorun Chen, Joel Jang, Yi Li, Siwei Han, Chaoqi Wang, Mingyu Ding, Dieter Fox, and Huaxiu Yao. Grape: Generalizing robot policy via preference alignment. *arXiv preprint arXiv:2411.19309*, 2024.

- Guanxing Lu, Wenkai Guo, Chubin Zhang, Yuheng Zhou, Haonan Jiang, Zifeng Gao, Yansong Tang, and Ziwei Wang. Vla-rl: Towards masterful and general robotic manipulation with scalable reinforcement learning. *arXiv preprint arXiv:2505.18719*, 2025a.
- Yuhui Chen, Shuai Tian, Shugao Liu, Yingting Zhou, Haoran Li, and Dongbin Zhao. Conrft: A reinforced fine-tuning method for vla models via consistency policy. *arXiv preprint arXiv:2502.05450*, 2025a.
- Murtaza Dalal, Tarun Chiruvolu, Devendra Chaplot, and Ruslan Salakhutdinov. Plan-seq-learn: Language model guided rl for solving long horizon robotics tasks. *arXiv preprint arXiv:2405.01534*, 2024.
- Jiarui Yang, Bin Zhu, Jingjing Chen, and Yu-Gang Jiang. Actor-critic for continuous action chunks: A reinforcement learning framework for long-horizon robotic manipulation with sparse reward. *arXiv preprint arXiv:2508.11143*, 2025.
- Adrià López Escoriza, Nicklas Hansen, Stone Tao, Tongzhou Mu, and Hao Su. Multi-stage manipulation with demonstration-augmented reward, policy, and world model learning. In *International Conference on Machine Learning*, 2025.
- Zixuan Chen, Jing Huo, Yangtao Chen, and Yang Gao. Robohorizon: An llm-assisted multi-view world model for long-horizon robotic manipulation, 2025b.
- Jiankai Sun, Aidan Curtis, Yang You, Yan Xu, Michael Koehle, Qianzhong Chen, Suning Huang, Leonidas Guibas, Sachin Chitta, Mac Schwager, and Hui Li. Arch: Hierarchical hybrid learning for long-horizon contact-rich robotic assembly, 2025.
- Qianzhong Chen, Justin Yu, Mac Schwager, Pieter Abbeel, Yide Shentu, and Philipp Wu. Sarm: Stage-aware reward modeling for long horizon robot manipulation. *arXiv preprint arXiv:2509.25358*, 2025c.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.
- Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, et al. A survey of reinforcement learning for large reasoning models. *arXiv preprint arXiv:2509.08827*, 2025a.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837, 2022.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- Hongyin Zhang, Zifeng Zhuang, Han Zhao, Pengxiang Ding, Hongchao Lu, and Donglin Wang. Reinbot: Amplifying robot visual-language manipulation with reinforcement learning. *arXiv preprint arXiv:2505.07395*, 2025b.

- Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-Jen Wang, Yucheng Hu, and Jianyu Chen. Improving vision-language-action model with online reinforcement learning. *arXiv preprint arXiv:2501.16664*, 2025b.
- Shuhan Tan, Kairan Dou, Yue Zhao, and Philipp Krähenbühl. Interactive post-training for vision-language-action models. *arXiv preprint arXiv:2505.17016*, 2025.
- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- Jijia Liu, Feng Gao, Bingwen Wei, Xinlei Chen, Qingmin Liao, Yi Wu, Chao Yu, and Yu Wang. What can rl bring to vla generalization? an empirical study. *arXiv preprint arXiv:2505.19789*, 2025.
- Junyang Shu, Zhiwei Lin, and Yongtao Wang. Rftf: Reinforcement fine-tuning for embodied agents with temporal feedback. *arXiv preprint arXiv:2505.19767*, 2025.
- Haozhan Li, Yuxin Zuo, Jiale Yu, Yuhao Zhang, Zhaohui Yang, Kaiyan Zhang, Xuekai Zhu, Yuchen Zhang, Tianxing Chen, Ganqu Cui, et al. Simplevla-rl: Scaling vla training via reinforcement learning. *arXiv preprint arXiv:2509.09674*, 2025a.
- Hongzhi Zang, Mingjie Wei, Si Xu, Yongji Wu, Zhen Guo, Yuanqing Wang, Hao Lin, Liangzhi Shi, Yuqing Xie, Zhexuan Xu, et al. Rlinf-vla: A unified and efficient framework for vla+ rl training. *arXiv preprint arXiv:2510.06710*, 2025.
- Kang Chen, Zhihao Liu, Tonghe Zhang, Zhen Guo, Si Xu, Hao Lin, Hongzhi Zang, Quanlu Zhang, Zhaofei Yu, Guoliang Fan, et al. π_{RL} : Online rl fine-tuning for flow-based vision-language-action models. *arXiv preprint arXiv:2510.25889*, 2025d.
- Physical Intelligence, Ali Amin, Raichelle Aniceto, Ashwin Balakrishna, Kevin Black, Ken Conley, Grace Connors, James Darpinian, Karan Dhabalia, Jared DiCarlo, et al. $\pi_{0.6}^*$: a vla that learns from experience. *arXiv preprint arXiv:2511.14759*, 2025.
- Yunfei Li, Xiao Ma, Jiafeng Xu, Yu Cui, Zhongren Cui, Zhigang Han, Liqun Huang, Tao Kong, Yuxiao Liu, Hao Niu, et al. Gr-rl: Going dexterous and precise for long-horizon robotic manipulation. *arXiv preprint arXiv:2512.01801*, 2025b.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- Dohyeong Kim, Hyeokjin Kwon, Junseok Kim, Gunmin Lee, and Songhwai Oh. Stage-wise reward shaping for acrobatic robots: A constrained multi-objective reinforcement learning approach. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025.
- Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, 1999.
- Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, et al. Evaluating real-world robot manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024a.
- Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander Hinrichsen, Xiaodi Yuan, Chen Bao, Xinsong Lin, Yulin Liu, Tse kai Chan, Yuan Gao, Xuanlin Li, Tongzhou Mu, Nan Xiao, Arnav Gurha, Viswesh Nagaswamy Rajesh, Yong Woo Choi, Yen-Ru Chen, Zhiao Huang, Roberto Calandra, Rui Chen, Shan Luo, and Hao Su. Maniskill3: Gpu parallelized robotics simulation and rendering for generalizable embodied ai. *Robotics: Science and Systems*, 2025.

- Sami Haddadin. The franka emika robot: A standard platform in robotics research [survey]. *IEEE Robotics & Automation Magazine*, 2024.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael S Ryoo, Grecia Salazar, Pannag R Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan H Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: Robotics Transformer for Real-World Control at Scale. In *Robotics: Science and Systems*, 2023.
- Xinghang Li, Peiyan Li, Minghuan Liu, Dong Wang, Jirong Liu, Bingyi Kang, Xiao Ma, Tao Kong, Hanbo Zhang, and Huaping Liu. Towards generalist robot policies: What matters in building vision-language-action models. *arXiv preprint arXiv:2412.14058*, 2024b.
- Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-action model. *arXiv preprint arXiv:2501.15830*, 2025.
- Zekun Qi, Wenyao Zhang, Yufei Ding, Runpei Dong, Xinqiang Yu, Jingwen Li, Lingyun Xu, Baoyu Li, Xialin He, Guofan Fan, Jiazhaoh Zhang, Jiawei He, Jiayuan Gu, Xin Jin, Kaisheng Ma, Zhizheng Zhang, He Wang, and Li Yi. Sofar: Language-grounded orientation bridges spatial reasoning and object manipulation. *Advances in neural information processing systems*, 2025.
- Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, Caroline Pascal, Martino Russi, Andres Marafioti, et al. Smolvla: A vision-language-action model for affordable and efficient robotics. *arXiv preprint arXiv:2506.01844*, 2025.
- Weifeng Lu, Minghao Ye, Zewei Ye, Ruihan Tao, Shuo Yang, and Bo Zhao. Robofac: A comprehensive framework for robotic failure analysis and correction. *arXiv preprint arXiv:2505.12224*, 2025b.

A. Training Settings and Evaluation Protocols

A.1. Supervised Fine-Tuning (SFT)

We initialize SFT models from two pretrained VLA backbones: OpenVLA-7B and the pi0.5_base. For each backbone, we optimize the action prediction objective using the AdamW optimizer with a learning rate of 1×10^{-5} and a constant schedule. SFT is trained for 100K steps per backbone on each benchmark, and the same training budget is used for all methods on the same backbone to ensure fair comparison. All image observations are resized to 224×224 before being fed into the VLA backbone.

A.2. Trajectory Preference Optimization: TPO and STA-TPO

For offline preference optimization, we start from SFT checkpoints of both OpenVLA-7B and pi0.5_base backbone, and construct datasets of paired successful (“chosen”) and failed (“rejected”) trajectories. Each trajectory consists of image observations, language instructions, and continuous robot actions, which are discretized into tokens using an action tokenizer built on top of the base vocabulary.

We use a DPO-style preference objective that increases the log-likelihood ratio between chosen and rejected trajectories relative to a frozen SFT reference model. This provides a stable and scalable way to perform trajectory-level preference alignment. Training uses AdamW with a learning rate of 2×10^{-5} , gradient accumulation, and 50K optimization steps per backbone to ensure fair comparison across baselines.

STA-TPO uses the same data pipeline, objective, and optimizer, but modifies the DPO logit difference by adding stage-aware margins computed from our stage calculator. This biases the preference updates toward stages where the chosen and rejected trajectories differ the most, while keeping easy stages largely unchanged, yielding a stage-aware variant that is directly comparable to standard TPO.

A.3. Reinforcement Learning: PPO and STA-PPO

We adopt an on-policy PPO framework. For SimplerEnv, we run 100 parallel environments with an episode horizon of 60 steps, yielding 6000 transitions per PPO update. Rollouts are stored in a separated replay buffer, and advantages are computed using generalized advantage estimation with a discount factor $\gamma = 0.99$ and GAE parameter $\lambda = 0.95$.

Optimization follows the standard clipped PPO objective. We use the AdamW optimizer with a policy learning rate of 1×10^{-4} , a value learning rate of 3×10^{-3} , and gradient accumulation over 20 steps to stabilize large-model training. Each update performs one PPO epoch over four minibatches, and the entropy coefficient is set to 0.0. Unless otherwise stated, ManiSkill3 uses the same optimization settings and training budget. All runs are trained for 600K environment steps to ensure fair comparison across baselines.

STA-PPO follows the same training pipeline and computational budget but augments the reward with dense stage-aware signals generated by the stage calculator. These structured shaping terms provide intermediate guidance that complements sparse task rewards, and the system additionally logs per-stage performance statistics during rollouts. Aside from this reward augmentation, STA-PPO is identical to PPO, enabling controlled comparisons of stage-aware credit assignment.

A.4. Evaluation Protocols

We evaluate all methods on both SimplerEnv-WidowX and ManiSkill3 benchmarks under a unified protocol. Each policy is tested over 300 evaluation episodes with deterministic decoding. A rollout

is considered successful if the environment-defined task condition is satisfied within the evaluation horizon (60 steps for `SimplerEnv` and 30 steps for `ManiSkill3`). Results are averaged over three random seeds. In addition to final success rates, we compute conditional stage success to diagnose at which manipulation stages policies succeed or fail.

A.5. Checkpoint Selection

We adopt a fixed-duration training schedule for all methods and use a consistent checkpoint-selection protocol across backbones and benchmarks.

For TPO and STA-TPO, models are trained for 50K optimization steps, and evaluated every 1,000 training steps on both `SimplerEnv` and `ManiSkill3`. The best-performing checkpoint under this evaluation protocol is reported.

For PPO and STA-PPO, models are trained for 600K steps. We run periodic evaluations every 6,000 environment steps on `SimplerEnv` and every 3,000 environment steps on `ManiSkill3`, selecting the checkpoint with the highest success rate.

For all ablation studies, we fix the total training duration and report the final checkpoint to ensure strict comparability across variants.

B. Experimental Setup and Implementation Details

B.1. Environments and Task Definitions

To evaluate robotic policies under diverse manipulation scenarios, we conduct experiments in two simulation environments: **`SimplerEnv-WidowX`** and **`ManiSkill3-Franka`**.

B.1.1. *SimplerEnv-WidowX*

We adopt the benchmark provided by `SimplerEnv` [Li et al., 2024a]. Within this environment, we test the following tasks:

- **Put the spoon on the towel** — The spoon is initialized at one corner of a 15×15 cm² square region on the tabletop; a towel is placed at the opposite corner. The spoon’s orientation alternates between horizontal and vertical across trials, requiring the gripper to adapt its pose accordingly. We perform 24 independent trials.
- **Put the carrot on the plate** — This task mirrors the previous one, with a carrot replacing the spoon and a plate replacing the towel. Similarly, 24 trials are executed.
- **Stack the green block on the yellow block** — A green block and a yellow block (both 3 cm cubes) are placed at different corners of a tabletop square. We evaluate two square sizes (side length 10 cm and 20 cm), yielding a total of 24 trials. This task tests precise grasping, lifting, and alignment to form a stable stack.
- **Put the eggplant into the yellow basket** — An eggplant is randomly placed in the right basin of a simulated sink, and a yellow basket is placed in the left basin. The eggplant’s position and orientation vary across trials (while avoiding basin edges to ensure graspability). A total of 24 trials are conducted. This task challenges the policy’s ability to deal with irregular shapes and variable object poses.

B.1.2. ManiSkill3-Franka

To further stress-test manipulation under contact-rich and long-horizon tasks, we use four tasks from ManiSkill3 [Tao et al., 2025] with a simulated Franka robot:

- **Stack Cube** — The robot must stack one cube on top of another, requiring accurate grasping, lifting, alignment, and stable placement.
- **Push Cube** — A cube is placed on the table, and the robot must push it toward a designated target region. This task emphasizes smooth trajectory execution and contact-rich control.
- **Pull Cube** — The robot grasps the cube on one side and pulls it (dragging) into the target region, demanding stable grasp maintenance under sliding contact.
- **Lift Peg Upright** — A peg is initialized lying flat on the table; the robot must grasp it, lift it, and reorient it so that it stands vertically upright. This task is challenging due to orientation constraints and the need for precise control to maintain balance.

B.2. Evaluation Protocol

For all experiments, we follow the evaluation procedure below:

- Each task is repeated across the full set of defined trials. For SimplerEnv we run 24 trials per task; for ManiSkill3 we follow the standard configurations defined by the benchmark.
- Success is determined based on task-specific completion criteria (e.g., correct placement, stable stack, upright peg, etc.).
- For tasks (or sequences) that involve multiple steps, we additionally record the trajectory length and measure efficiency, robustness, and consistency across trials.
- Final results (success rates, trajectory lengths) are reported as mean \pm standard deviation over all trials, ensuring statistical reliability.

B.3. Implementation Details

- Both SimplerEnv and ManiSkill3 are used with their default simulation configurations. We do not modify the underlying environment logic — only vary object initial positions and orientations as described.
- For each trial, object poses (positions, orientations when applicable) are randomized within the constraints defined for each task to ensure diversity across trials.
- Policies operate in an end-to-end fashion (observations \rightarrow actions), without manual heuristics or task-specific engineering.
- All reported metrics are averaged over the full set of trials. When appropriate, we also report standard deviations to reflect variability.

B.4. Data Collection

Our data collection process, which is similar to the methodology used in the main paper, is designed to generate high-quality data for both supervised and preference-based learning. All data is collected specifically for the selected ManiSkill3 tasks.

- **Expert Demonstrations:** For each task, we generate **100 high-quality demonstration trajectories**. These trajectories are produced using the MPLib motion planner, ensuring kinematically feasible and efficient paths to task completion. Following the findings of the main paper, we

apply an action filtering technique to this data, removing idle actions where the end-effector pose changes by a negligible amount. This preprocessing step is crucial for mitigating the issue of trained SFT policies getting stuck during execution.

- **Preference Pairs:** For methods requiring preference data (e.g., TPO), we generate **50 trajectory preference pairs** per task. In the case of SimplerEnv, trajectories are sampled for each task using the Octo model. These pairs are obtained by sampling two trajectories from the Octo-collected dataset, and assigning preference labels based on cumulative rewards (e.g., successful completion vs. failure).

C. Stage-wise Cost and Potential Definitions

This appendix provides the complete definitions of the stage-wise costs $\ell_k(\tau^{(k)})$ and progress potentials $\Phi_k(s)$ used in STARE. Each stage corresponds to a geometric manipulation primitive. Stage costs—used in STA-TPO—evaluate how well a trajectory segment achieves the geometric goal of a stage, while stage potentials—used in STA-PPO—provide dense, per-step shaping signals that reflect normalized progress within a stage. All potentials follow a normalized sigmoid form,

$$\sigma(z) = \frac{1}{1 + e^{-z}},$$

and lie in $[0, 1]$. For stages involving orientation, we additionally use a normalized rotational alignment measure $g(R_1, R_2) \in [0, 1]$, where $g(R_1, R_2)$ denotes a normalized orientation alignment measure.

Stage Decomposition Across Tasks. Different manipulation tasks instantiate different subsets of these generic manipulation primitives. We use a minimal, task-agnostic decomposition based on geometric event boundaries that can be reliably detected from state observations:

- **Pick-Place:** Reach \rightarrow Grasp \rightarrow Transport \rightarrow Place.
- **Push:** Reach \rightarrow Push \rightarrow Goal.
- **Pull:** Reach \rightarrow Pull \rightarrow Goal.
- **Peg Upright:** Reach \rightarrow Grasp \rightarrow Lift \rightarrow Upright.

Each stage is associated with a well-defined geometric objective (e.g., approach, alignment, elevation, planar displacement, fine-grained goal adjustment). The following sections list the complete cost and potential functions for all stages used in STARE.

C.1. Reach

Goal. Guide the end-effector toward the target object until it enters the vicinity where grasp alignment becomes feasible. This stage captures the coarse approach motion before any contact or fine-grained adjustment occurs.

Cost. We quantify reaching quality using the Euclidean distance between the end-effector position and the object’s grasp point:

$$d_{\text{reach}}(t) = \|x^{\text{ee}}(t) - x^{\text{obj}}(t)\|.$$

The stage cost is the average approach error over the stage segment:

$$\ell_{\text{Reach}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{reach}}(t).$$

Potential. To provide dense shaping during the approach, we define a normalized potential that increases as the end-effector draws closer to the object:

$$\Phi_{\text{Reach}}(s) = \sigma\left(1 - \frac{d_{\text{reach}}(s)}{d_{\text{reach}}^{\max}}\right).$$

Here d_{reach}^{\max} is a geometry-derived normalization scale, chosen as the object’s characteristic size:

$$d_{\text{reach}}^{\max} = L_{\text{obj}},$$

where L_{obj} denotes the object’s side length (e.g., cube size in ManiSkill3) or, more generally, its bounding-box diameter. The same d_{reach}^{\max} is used as the threshold for the Reach \rightarrow Grasp transition, ensuring that approach shaping and stage segmentation operate on a consistent geometric scale.

C.2. Grasp

Goal. Establish and maintain a stable grasp. In ManiSkill environments, the end-effector approaches the object, aligns with the grasp point, and closes the gripper to securely hold the object.

Cost. The Grasp stage evaluates how well the end-effector pose matches the intended grasp configuration. We measure this using the distance between the tool-center point (TCP) and the object’s grasp point:

$$d_{\text{pose}}(t) = \left\|x^{\text{tcp}}(t) - x^{\text{obj}}(t)\right\|.$$

The stage cost is the average geometric misalignment:

$$\ell_{\text{Grasp}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{pose}}(t).$$

Potential. We define a normalized potential that increases as the TCP approaches the grasp point:

$$\Phi_{\text{Grasp}}(s) = \sigma\left(1 - \frac{d_{\text{pose}}(s)}{d_{\text{grasp}}}\right),$$

where $\sigma(\cdot)$ is the logistic sigmoid and d_{grasp} is a task-agnostic normalization scale determined by object geometry:

$$d_{\text{grasp}} = L_{\text{obj}},$$

with L_{obj} the characteristic object size (e.g., the cube side length in ManiSkill3). This scale is also used as the threshold for the Reach \rightarrow Grasp stage transition, ensuring consistent geometry-aware shaping without per-task tuning.

C.3. Transport

Goal. Move the grasped object toward its target goal pose while maintaining a stable grasp. This stage captures the coarse relocation of the object once it has been lifted or secured.

Cost. Transport quality is measured using the distance between the object’s current position and the goal location:

$$d_{\text{trans}}(t) = \left\| x^{\text{obj}}(t) - x^{\text{goal}} \right\|.$$

The stage cost averages this residual goal error:

$$\ell_{\text{Transport}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{trans}}(t).$$

Potential. To provide dense shaping toward the goal, we define:

$$\Phi_{\text{Transport}}(s) = \sigma \left(1 - \frac{d_{\text{trans}}(s)}{d_{\text{trans}}^{\text{max}}} \right).$$

The normalization scale $d_{\text{trans}}^{\text{max}}$ is set to the characteristic task displacement:

$$d_{\text{trans}}^{\text{max}} = \left\| x_{\text{init}}^{\text{obj}} - x^{\text{goal}} \right\|,$$

i.e., the distance between the object’s initial position and its goal. This provides a task-agnostic scale reflecting the required transport distance, and aligns potential shaping with the physical extent of the manipulation.

C.4. Place

Goal. Position the object precisely at the goal location and stabilize it. This stage captures the fine-grained alignment after coarse transport has moved the object near its target.

Cost. Placement quality is measured by the residual distance between the object’s current position and the goal:

$$d_{\text{place}}(t) = \left\| x^{\text{obj}}(t) - x^{\text{goal}} \right\|.$$

The stage cost averages this near-goal deviation:

$$\ell_{\text{Place}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{place}}(t).$$

Potential. We define a normalized potential encouraging precise placement:

$$\Phi_{\text{Place}}(s) = \sigma \left(1 - \frac{d_{\text{place}}(s)}{d_{\text{place}}^{\text{max}}} \right).$$

The normalization scale is chosen as the object’s characteristic size:

$$d_{\text{place}}^{\text{max}} = L_{\text{obj}},$$

which reflects the resolution required for fine positioning and provides a task-agnostic geometric scale.

C.5. Push & Pull

Goal. Move the object toward the target goal while maintaining stable contact with the end-effector. This stage unifies push- and pull-based planar manipulation, as both correspond to contact-induced object translation in the plane.

Cost. Progress is measured using the residual distance between the object and its target:

$$d_{pp}(t) = \left\| x^{\text{obj}}(t) - x^{\text{goal}} \right\|.$$

The stage cost averages this residual error:

$$\ell_{\text{Push\&Pull}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{pp}(t).$$

Potential. We define a normalized potential that increases as the object approaches the goal:

$$\Phi_{\text{Push\&Pull}}(s) = \sigma \left(1 - \frac{d_{pp}(s)}{d_{pp}^{\text{max}}} \right),$$

where the normalization scale is chosen as the required planar displacement:

$$d_{pp}^{\text{max}} = \left\| x_{\text{init}}^{\text{obj}} - x^{\text{goal}} \right\|.$$

This yields a task-agnostic geometric scale for shaping and is consistent with the Transport stage, differing only by the presence of contact.

C.6. Lift

Goal. Lift the object from the table to a target height z_{goal} , ensuring that it is safely elevated above obstacles for subsequent manipulation.

Cost. We define the lift deviation as the vertical residual:

$$d_{\text{lift}}(t) = |z_{\text{obj}}(t) - z_{\text{goal}}|.$$

The stage cost averages this error:

$$\ell_{\text{Lift}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{lift}}(t).$$

Potential. We define a normalized shaping potential:

$$\Phi_{\text{Lift}}(s) = \sigma \left(1 - \frac{|z_{\text{obj}}(s) - z_{\text{goal}}|}{d_{\text{lift}}^{\text{max}}} \right),$$

where

$$d_{\text{lift}}^{\text{max}} = z_{\text{goal}} - z_{\text{table}}$$

is the required lifting displacement. This provides a task-agnostic geometric scale and yields dense shaping that increases smoothly as the object approaches its target height.

C.7. Upright

Goal. Align the object’s orientation with an upright target orientation. This stage captures fine-grained rotational adjustment after the object has been lifted or placed near its desired pose.

Cost. We measure uprightness using the angular deviation between the object orientation $R_{\text{obj}}(t)$ and the target upright orientation R_{upright} . Let

$$d_{\text{upright}}(t) = \arccos\left(\frac{\text{tr}(R_{\text{upright}}^\top R_{\text{obj}}(t)) - 1}{2}\right),$$

which equals the geodesic rotation distance on $\text{SO}(3)$. The stage cost averages this orientation error:

$$\ell_{\text{Upright}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{upright}}(t).$$

Potential. We define a normalized shaping potential that increases as the object approaches its upright orientation:

$$\Phi_{\text{Upright}}(s) = \sigma\left(1 - \frac{d_{\text{upright}}(s)}{d_{\text{upright}}^{\max}}\right).$$

The normalization scale $d_{\text{upright}}^{\max}$ corresponds to the maximum possible rotational deviation:

$$d_{\text{upright}}^{\max} = \pi,$$

which is the geodesic diameter of $\text{SO}(3)$. This choice yields a task-agnostic, geometry-grounded scale that applies to any upright-orientation task.

C.8. Goal (Push/Pull)

Goal. Ensure that the object reaches the target goal region and remains stably within it.

Cost. We measure goal attainment using the residual distance between the object and the goal:

$$d_{\text{goal}}(t) = \|x^{\text{obj}}(t) - x^{\text{goal}}\|.$$

The stage cost averages this error:

$$\ell_{\text{Goal}}(\tau^{(k)}) = \frac{1}{T_k} \sum_{t \in \tau^{(k)}} d_{\text{goal}}(t).$$

Potential. A normalized potential provides dense shaping near the goal:

$$\Phi_{\text{Goal}}(s) = \sigma\left(1 - \frac{d_{\text{goal}}(s)}{d_{\text{goal}}^{\max}}\right),$$

where

$$d_{\text{goal}}^{\max} = L_{\text{obj}},$$

the characteristic object size. This creates a fine-scale potential landscape around the goal region, allowing smooth convergence without relying on a binary indicator.

C.9. Notation Summary

Table 3 lists all variables used in the stage-wise cost and potential definitions across all tasks.

Symbol	Meaning
$x^{ee}(t)$	End-effector (TCP) position at time t .
$x^{obj}(t)$	Object center or grasp-point position.
x^{goal}	Target goal position of the object.
x_{init}^{obj}	Object position at episode start.
$d_{reach}(t)$	TCP–object distance during Reach.
$d_{pose}(t)$	TCP–object alignment error during Grasp.
$d_{trans}(t)$	Object–goal residual distance for Transport.
$d_{place}(t)$	Near-goal deviation during Place.
$d_{pp}(t)$	Residual distance for Push&Pull.
$d_{lift}(t)$	Vertical height error.
$d_{upright}(t)$	Geodesic rotation distance on $SO(3)$.
d^{max}	Normalization scales (stage-dependent).
L_{obj}	Characteristic object size (e.g., cube side length).
T_k	Number of steps in stage $\tau^{(k)}$.
$R_{obj}(t)$	Object orientation at time t .
$R_{upright}$	Target upright orientation.

Table 3 | Summary of notation used across all stage definitions.

D. Broader Impacts

Our framework aims to improve safety and reliability of robot manipulation by aligning policies with human-preferred, semantically correct behaviors. Potential risks include overfitting to biased preferences and misuse in unsafe settings; we discuss mitigations in the main text and Appendix.