

# RLinf-VLA: A Unified and Efficient Framework for Reinforcement Learning of Vision-Language-Action Models

Anonymous CVPR submission

Paper ID 4

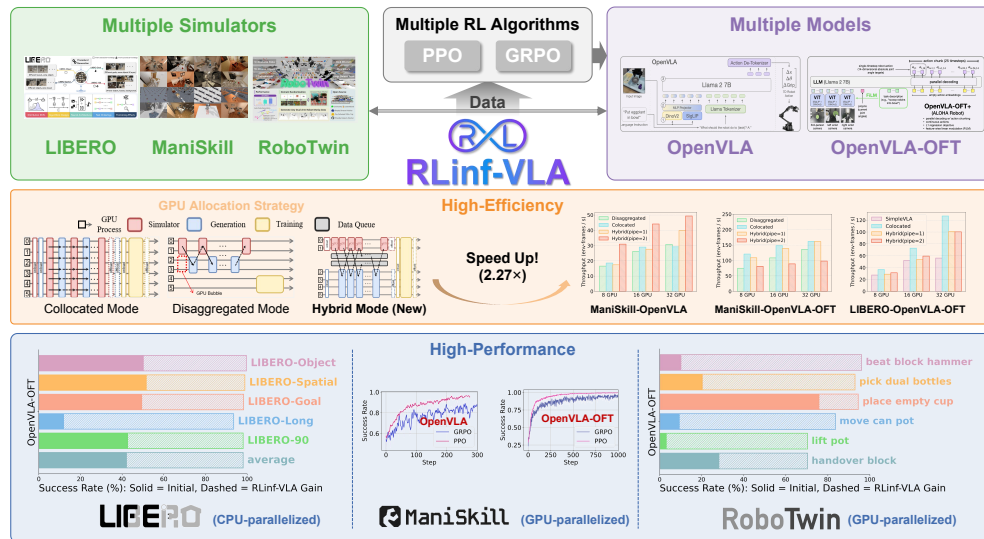


Figure 1. Overview of RLinf-VLA. Built on a unified interface, RLinf-VLA seamlessly supports diverse VLA architectures, multiple RL algorithms, and various simulators. It provides three GPU allocation modes: *collocated*, *disaggregated*, and a novel *hybrid* mode, speeding up training by  $2.27\times$  compared to the baseline. A single unified model achieves 98.11% success on 130 LIBERO tasks and 97.66% on 25 ManiSkill tasks. RLinf-VLA achieves an average 84.63% success on 6 RoboTwin tasks.

## Abstract

001 Recent advances in vision-language-action (VLA) models  
 002 have motivated the extension of their capabilities to em-  
 003 bodied settings, where reinforcement learning (RL) offers  
 004 a principled way to optimize task success through interac-  
 005 tion. However, existing methods remain fragmented, lack-  
 006 ing both a unified platform for fair comparison across ar-  
 007 chitectures and algorithms and an efficient system design  
 008 for scalable training. To address these challenges, we in-  
 009 troduce RLinf-VLA, a unified and efficient framework for  
 010 scalable RL training of VLA models. RLinf-VLA achieves  
 011 unification by providing a unified interface that standar-  
 012 dizes the integration of diverse VLA architectures, multiple  
 013 RL algorithms, and heterogeneous simulators, enabling ex-  
 014 tensibility. To ensure efficiency, the system adopts a flex-  
 015 ible resource allocation architecture for rendering, infer-

ence, and training workloads in RL pipelines. In particular,  
 for GPU-parallelized simulators, RLinf-VLA introduces a  
 hybrid fine-grained pipeline allocation strategy, yielding a  
 $1.61\times$ – $1.88\times$  training speedup. Using this unified system,  
 models trained with RLinf-VLA demonstrate consistent per-  
 formance improvements of approximately 20–85% across  
 multiple simulation benchmarks, including LIBERO, Mani-  
 Skill, and RoboTwin. Furthermore, we distill a set of train-  
 ing practices for effective RL-based VLA training. We posi-  
 tion RLinf-VLA as a foundational system to enable efficient,  
 unified, and reproducible research in embodied intelligence.

016  
 017  
 018  
 019  
 020  
 021  
 022  
 023  
 024  
 025  
 026  
 027

## 1. Introduction

Vision-Language-Action (VLA) models represent a new  
 generation of foundation models that aim to unify per-

028  
 029  
 030

031 ception, language understanding, and embodied control  
032 [9, 25]. By leveraging vision-language models pretrained  
033 on internet-scale data and further training them on large,  
034 heterogeneous robot demonstration datasets [16, 30], VLAs  
035 can map raw sensor observations and natural language in-  
036 structions directly to robot actions. This paradigm has  
037 demonstrated strong generalization across a wide range of  
038 tasks [2, 15, 17, 23, 34, 40, 41].

039 However, deploying VLA models typically requires  
040 post-training to mitigate distribution mismatch between  
041 training data and deployment environments. Without ef-  
042 fective post-training, performance can degrade significantly  
043 when models encounter novel states, instructions, or execu-  
044 tion conditions.

045 Reinforcement learning (RL) has emerged as an increas-  
046 ingly important post-training paradigm, as it directly opti-  
047 mizes cumulative task rewards through trial-and-error in-  
048 teraction [37]. In contrast to supervised fine-tuning (SFT),  
049 another commonly used post-training approach, RL en-  
050 ables exploration beyond narrow expert demonstrations and  
051 equips policies with corrective and adaptive behaviors. Re-  
052 cent studies indicate that RL fine-tuning can yield stronger  
053 out-of-distribution generalization than SFT, particularly in  
054 terms of semantic alignment and execution robustness [22].

055 Despite its promise, applying RL to VLA models re-  
056 mains largely small-scale or fragmented [22, 24], mak-  
057 ing the development of new RL algorithms for VLAs ex-  
058 pensive and cumbersome. Although SimpleVLA-RL [20]  
059 enables large-scale RL training for VLAs by building on  
060 VeRL [36], an RL codebase originally designed for large  
061 language models, it lacks system-level optimizations tai-  
062 lored to embodied settings, where simulators compete with  
063 model inference and learning for GPU resources. This  
064 limitation constrains scalability and substantially increases  
065 training time.

066 Moreover, online RL requires repeated and tightly cou-  
067 pled model–environment interactions. Without a system de-  
068 sign explicitly tailored to this execution pattern, significant  
069 GPU idle time and pipeline bubbles can arise, leading to in-  
070 efficient resource utilization. Finally, existing works often  
071 rely on disparate model architectures, RL algorithms, and  
072 evaluation protocols, making fair comparison difficult and  
073 hindering the extraction of general principles for RL-based  
074 VLA training.

075 To address these challenges, we present **RLinf-VLA**, a  
076 *unified* and *efficient* framework for scalable reinforcement  
077 learning of vision-language-action models. Our main con-  
078 tributions are summarized as follows:

079 • **Unified system abstraction.** RLinf-VLA provides a  
080 unified interface that supports multiple robotic simula-  
081 tors (ManiSkill [39], LIBERO [21], RoboTwin [8]), di-  
082 verse VLA architectures (OpenVLA [17], OpenVLA-  
083 OFT [18]), and reinforcement learning algorithms

(PPO [33], GRPO [35]). The system exposes three ex- 084  
085 ecution modes, including a novel hybrid GPU allocation  
086 mode, enabling scalable and configurable training across  
087 heterogeneous setups.

- 088 • **Efficient system and algorithm design.** RLinf-VLA  
089 introduces hybrid fine-grained pipelining for GPU-  
090 parallelized simulators and collocated execution for CPU-  
091 parallelized simulators, substantially improving training  
092 throughput, with speedups of up to 2.27×. In addition, we  
093 incorporate a set of algorithmic optimizations that further  
094 enhance training efficiency and stability.
- 095 • **Strong empirical performance and generalization.**  
096 Using a single unified model, RLinf-VLA achieves a  
097 98.11% success rate on 130 LIBERO tasks and 97.66%  
098 on 25 ManiSkill tasks. RLinf-VLA also achieves an aver-  
099 age success rate of 84.63% on six RoboTwin tasks, real-  
100 izing an average performance improvement of 63.75% on  
101 these tasks, which demonstrates its powerful performance  
102 and generalization capabilities in post-training.
- 103 • **Open and extensible platform.** RLinf-VLA is released  
104 as an open-source and actively maintained platform, pro-  
105 viding a practical foundation to accelerate, standardize,  
106 and scale reinforcement learning research for embodied  
107 intelligence.

## 108 2. Related Work

### 109 2.1. Reinforcement Learning for VLA

110 Vision-Language-Action models (VLAs) [1, 2, 4–6, 15, 17,  
111 18, 43, 46] represent a multimodal foundation model archi-  
112 tecture designed for robotics. RL is increasingly adopted  
113 to address the limitations of static imitation learning in  
114 VLAs. GRAPE [44] employs Direct Preference Optimiza-  
115 tion (DPO [31]) on partial trajectories to align models with  
116 human intent. Iterative approaches [12, 38] combine SFT  
117 with RL stages to balance stability and performance. VLA-  
118 RL [24] utilizes Proximal Policy Optimization (PPO [33])  
119 to exhibit significantly stronger generalization to unseen  
120 objects and environments compared to their SFT counter-  
121 parts [22]. SimpleVLA-RL [20] and related variants [38]  
122 apply Group Relative Policy Optimization (GRPO [35]) to  
123 improve training efficiency and performance without com-  
124 plex reward design.  $\pi_{RL}$  [7] implements the policy gra-  
125 dient algorithm to flow-matching model for the first time.  
126  $\pi_{0.6}^*$  [14] introduces RECAP to train a binarized value func-  
127 tion and VLA. A unified and efficient infrastructure is es-  
128 sential for advancing VLA via RL, yet current implemen-  
129 tations remain fragmented [22, 24]. Although frameworks  
130 such as SimpleVLA-RL [20] have emerged, they largely  
131 follow generic LLM RL post-training paradigms and fail to  
132 adequately account for the system-level optimizations and  
133 interface extensibility required by embodied intelligence  
134 tasks.

135 **2.2. Simulators**

136 Simulators [3, 10, 19, 26–29] offer scalable, safe, and con-  
 137 trollable environments that bypass physical hardware con-  
 138 straints. They facilitate massive parallel data generation and  
 139 automated resets, essential for training and evaluating gen-  
 140 eralizable VLA policies. Notable platforms include Man-  
 141 iSkill [39] for physics-rich manipulation, LIBERO [21] for  
 142 instruction-conditioned reasoning, and RoboTwin [8] for bi-  
 143 manual tasks with domain randomization. Consequently,  
 144 the construction and utilization of simulators are indispens-  
 145 able for RL. However, inconsistent interfaces across simu-  
 146 lators often necessitate extensive code restructuring. A uni-  
 147 fied interface is therefore critical to enable seamless switch-  
 148 ing between diverse simulation benchmarks.

149 **3. Preliminary**150 **3.1. Reinforcement Learning**

151 To guide the design of our efficient training infrastructure,  
 152 we first formalize the problem of finetuning VLA mod-  
 153 els using Reinforcement Learning (RL). We model vision-  
 154 based manipulation tasks as a partially observable Markov  
 155 decision process (POMDP). A POMDP extends the stan-  
 156 dard MDP and is defined by the tuple  $(\mathcal{S}, \mathcal{A}, P, r, \gamma, \Omega, O)$ ,  
 157 where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  the action space (dis-  
 158 crete or continuous, potentially chunked),  $P(s' | s, a)$   
 159 the state transition function,  $r(s, a)$  the reward function,  
 160  $\gamma \in [0, 1]$  the discount factor,  $\Omega$  the observation space, and  
 161  $O(o | s)$  the observation function mapping a state  $s$  to an  
 162 observation  $o$ .

163 At each timestep  $t$ , the agent samples an action  $a_t \sim$   
 164  $\pi_\theta(\cdot | o_t)$  based on the current observation  $o_t \sim O(\cdot | s_t)$ .  
 165 The environment then executes  $a_t$ , transitions to the next  
 166 state  $s_{t+1}$  via  $P$ , and produces the next observation  $o_{t+1}$ .  
 167 The goal of RL is to optimize the parameterized policy  $\pi_\theta$  so  
 168 as to maximize the expected cumulative discounted reward  
 169 over trajectories  $\tau = (s_0, a_0, \dots, s_T)$ :

$$170 \quad J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right]. \quad (1)$$

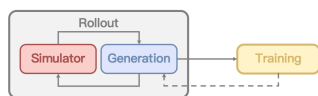
171 **3.2. Pipeline of Reinforcement Learning for VLA**  
172 **models**

Figure 2. Pipeline of Reinforcement Learning for VLA models

173 The RL pipeline for VLA models comprises three dis-  
 174 tinct components: *Generation*, *Simulator*, and *Training*, as

175 illustrated in Figure 2. GPU resources are primarily al-  
 176 located between the rollout phase (*Generation* and *Simu-  
 177 lator*) and the optimization phase (*Training*). A key dis-  
 178 tinction from LLM-based RL lies in the hierarchical action  
 179 formulation. While LLMs typically operate on single to-  
 180 kens [35], VLA policies often predict a single action [17]  
 181 or a *chunk* of actions to ensure smooth control [45]. In  
 182 our framework, we unify these representations by treating  
 183 a single action as a chunk of size one. Crucially, this in-  
 184 troduces a three-level hierarchy: *Chunk*  $\rightarrow$  *Atomic Action*  
 185  $\rightarrow$  *Token*. Each atomic action comprises multiple tokens,  
 186 where each token typically corresponds to a specific dimen-  
 187 sion of the robot’s action space (e.g., end-effector pose or  
 188 joint angles). Moreover, the integration of simulators im-  
 189 poses substantial computational overhead on the generation  
 190 process, demanding significant CPU and GPU resources for  
 191 physics simulation and rendering. The interaction typically  
 192 proceeds as follows [18]: the *Generation* module infers an  
 193 action chunk based on the current observation; the *Simula-  
 194 tor* executes this chunk and returns the observation from the  
 195 final timestep. This loop continues until trajectory collec-  
 196 tion is complete, after which the *Training* module updates  
 197 the policy using the gathered data.

198 **4. Method**

199 Section 4.1 introduces flexible GPU Allocation Strategies  
 200 for improved resource scheduling across different simulator  
 201 types. Section 4.2 presents a Unified Interface that seam-  
 202 lessly integrates diverse simulators, VLA models, and RL  
 203 algorithms. Building on this system design, Section 4.3 dis-  
 204 tills key Design Choices for scalable RL training of VLA  
 205 models.

206 **4.1. GPU Allocation Strategies**

207 The RL training pipeline (described in Section 3.2) for VLA  
 208 involves varying resource demands from *Training*, *Gener-  
 209 ation*, and *Simulator*. Resource bottlenecks depend heav-  
 210 ily on the simulator type: CPU-parallelized simulators are  
 211 typically CPU-bound, using GPUs primarily for rendering  
 212 and inference, whereas GPU-parallelized simulators exe-  
 213 cute simulation, rendering, and inference entirely on the  
 214 GPU. While the latter offers higher throughput, it creates  
 215 severe contention for GPU memory and compute. To max-  
 216 imize efficient utilization across these diverse setups, flex-  
 217 ible and optimized GPU allocation strategies are essential.  
 218 Our framework supports flexible and easily configurable al-  
 219 location modes: *collocated*, *disaggregated*, and a novel *hy-  
 220 brid* mode.

221 **4.1.1. Collocated Allocation**

222 In this mode (Figure 3), all components co-exist on the same  
 223 set of GPUs. The original version of collocated mode aimed  
 224 to maximize GPU utilization by maintaining only one com-

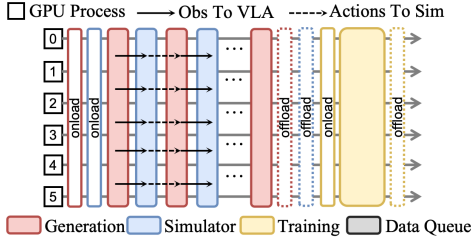


Figure 3. Collocated Mode.

225 ponent on all GPUs at any given time. When a component  
 226 finished its computation, it would be offloaded from GPU  
 227 to CPU memory, and then onloaded again when needed for  
 228 its next task. However, under the embodiment setting, both  
 229 *Simulator* and *Generation* require GPU resources and inter-  
 230 act frequently in an iterative manner. The overhead intro-  
 231 duced by repeated offload-onload operations at each inter-  
 232 action becomes prohibitively large and unacceptable.

233 Therefore, we implement a modified collocated strategy  
 234 where offload and onload operations for the *Simulator* and  
 235 *Generation* occur only at the beginning and end of the roll-  
 236 out phase, as shown in the figure. While this approach suc-  
 237 cessfully avoids the frequent offload-onload overhead dur-  
 238 ing rollout, it cannot fully utilize GPU resources and re-  
 239 mains sub-optimal. Since *Generation* and *Simulator* must  
 240 interact iteratively, they spend significant time waiting for  
 241 each other, occupying GPU memory without performing  
 242 computation, leading to resource wastage and limited scal-  
 243 ability.

#### 244 4.1.2. Disaggregated Allocation

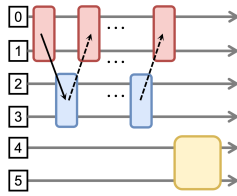


Figure 4. Disaggregated Mode (see Figure 3 for legend).

245 In this mode, each component is assigned to a distinct  
 246 (possibly multi-GPU) partition, with no overlap across  
 247 components. Figure 4 depicts the central idea. This en-  
 248 sures that every component can fully exploit its allocated  
 249 resources. This mode is simple to implement, but it lead  
 250 to GPU underutilization due to inter-component dependen-  
 251 cies. For example, in Figure 4, the gpu 4, 5 are completely  
 252 idle during the rollout phase util training.

#### 253 4.1.3. Hybrid Allocation with Fine-grained Pipelining

254 To overcome the drawbacks of the above modes, we pro-  
 255 pose a hybrid strategy: *hybrid allocation* combined with

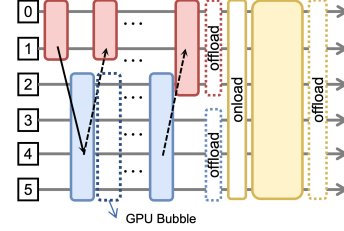


Figure 5. Hybrid Mode (see Figure 3 for legend).

256 *fine-grained pipelining*. In hybrid allocation, components  
 257 can flexibly select GPUs. A typical configuration is to as-  
 258 sign *Generation* and *Simulator* to different GPU partitions,  
 259 while allowing *Training* to utilize all GPUs. Figure 5 shows  
 260 an illustration. However, the resources remain underuti-  
 261 lized. For instance, the *Simulator* must wait for actions  
 262 generated by the *Generation*, leaving some GPUs idle and  
 creating “GPU bubbles”.

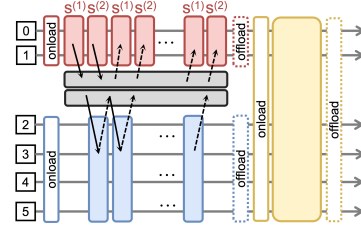


Figure 6. Hybrid Mode with Fine-grained Pipelining (see Figure 3 for legend).

263 On top of this, we introduce fine-grained pipelining  
 264 to mitigate bubbles caused by inter-component depen-  
 265 dencies. Specifically, a simulator instance on one GPU  
 266 is partitioned into multiple sub-simulators, denoted as  
 267  $S^{(1)}, S^{(2)}, \dots, S^{(k)}$ . The pipeline proceeds as follows,  
 268 where superscripts denote simulator indices and subscripts  
 269 indicate timesteps:  
 270

- 271 1. At step  $t = 0$ ,  $S^{(1)}$  generates the initial observation  $o_0^{(1)}$ ,  
 272 which is sent to the *Generation* component to produce  
 273 action  $a_0^{(1)}$ .
- 274 2. Meanwhile,  $S^{(2)}$  generates  $o_0^{(2)}$  in parallel.
- 275 3. Once  $a_0^{(1)}$  is ready, it is fed back into  $S^{(1)}$  to produce  
 276 the next observation  $o_1^{(1)}$ , while  $o_0^{(2)}$  is simultaneously  
 277 processed by the actor to generate  $a_0^{(2)}$ .

278 This scheduling allows *Simulator* and *Generation* to run  
 279 concurrently, reducing idle time while preserving correct-  
 280 ness. In this way, our framework avoids the frequent of-  
 281 floading overhead of collocated allocation while also elimi-  
 282 nating the GPU bubbles that occur in disaggregated alloca-  
 283 tion. Figure 5 provides a schematic illustration of the hybrid  
 284 allocation mode with fine-grained pipelining ( $k = 2$ ).

## 285 4.2. A Unified Interface

286 To conveniently reuse existing simulators, models, and al-  
287 gorithms, RLInfff-VLA provides a unified interface de-  
288 signed for high extensibility and scalability.

### 289 4.2.1. Multiple Simulators Support.

290 RLInff-VLA supports multiple robotic simulators, includ-  
291 ing ManiSkill [39], LIBERO [21], and RoboTwin [8]. To fa-  
292 cilitate this, we provide a unified interface that standardizes  
293 environment interactions across different backends, such as  
294 ManiSkill and RoboTwin’s GPU-parallelized environments  
295 and LIBERO’s CPU-parallelized wrappers. This interface  
296 includes standard Gym-style APIs with built-in support for  
297 action chunking and flexible episode termination. Further-  
298 more, we define a unified interface for observation and ac-  
299 tion to ensure seamless compatibility with various algo-  
300 rithms and models.

### 301 4.2.2. Multiple Models Support

302 Benefiting from a unified data processing pipeline across  
303 simulators, our framework facilitates the seamless integra-  
304 tion of existing models by requiring only a standardized in-  
305 teraction interface. We support diverse VLA architectures,  
306 exemplified in this study by OpenVLA [17] and OpenVLA-  
307 OFT [18], which represent configurations with single-step  
308 and multi-step action chunking, respectively. To enable  
309 efficient fine-tuning, we integrate Low-Rank Adaptation  
310 (LoRA) [13].

### 311 4.2.3. Multiple Algorithms Support

312 Our framework provides support for multiple reinforcement  
313 learning algorithms, with an initial focus on Proximal Pol-  
314 icy Optimization (PPO) [33] and Group Relative Policy Op-  
315 timization (GRPO) [11, 35]. The core functions for on-  
316 policy RL algorithms are the advantage function and the  
317 loss function. In our framework, we can configure the al-  
318 gorithms by specifying these two functions only, without  
319 redundant code.

## 320 4.3. Design Choices for Algorithms

321 RL algorithms for large models involve more complex de-  
322 sign choices than those for smaller models. For VLAs, we  
323 adopt methodologies from both LLMs and robotics. In Sec-  
324 tion 4.3.1, we introduce general features applicable to all al-  
325 gorithms. Subsequently, in Section 4.3.2 and Section 4.3.3,  
326 we detail specific design choices for PPO and GRPO, re-  
327 spectively. Consistent with our codebase style, all these  
328 features can be easily configured by changing values in the  
329 configuration file.

### 330 4.3.1. Multi-Granularity Calculation Support

331 Our framework supports advantage estimation and log-  
332 probability computation at multiple levels of granularity,

enabling seamless adaptation to different algorithmic re-  
quirements.

For advantage estimation, we provide two formula-  
tions for assigning advantages to action chunks  $c_t =$   
 $(a_{t,1}, \dots, a_{t,C})^1$ . In the chunk-level formulation, the entire  
chunk is treated as a single macro-action and is assigned a  
summed reward and a shared advantage. In contrast, the  
action-level formulation assigns individual rewards and ad-  
vantages to each atomic action  $a_{t,i}$  within the chunk.

Similarly, for log-probability computation, we support  
three corresponding granularities:

$$\text{Chunk-level: } \pi(c_t|o_t) = \prod_{i=1}^C \pi(a_{t,i}|o_t, a_{t:i-1}), \quad 344$$

$$\text{Action-level: } \pi(a_{t,i}|o_t, a_{t:i-1}) = \prod_{j=1}^M \pi(d_{t,i,j}|o_t, d_{t,i:j-1}), \quad 345$$

$$\text{Token-level: } \pi(d_{t,i,j}|o_t, d_{t,i:j-1}), \quad 346$$

where  $d_{t,i,j}$  denotes the  $j$ -th token of the  $i$ -th atomic action. 347

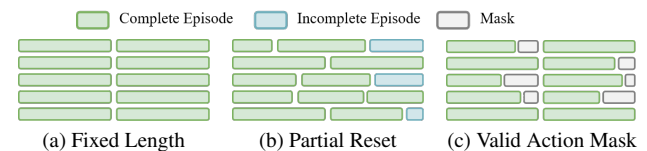


Figure 7. Illustration of the three rollout modes.

### 438 4.3.2. Design Choices for PPO

Scaling PPO to VLA models introduces several design chal-  
lenges. We summarize the key design decisions adopted in  
our framework below.

**Partial Reset Support** During rollouts, a sub-  
environment typically terminates either upon reaching the  
maximum episode length or upon successfully completing a  
task. We consider two standard strategies for handling such  
terminations. In the *Fixed Episode Length* mode, all sub-  
environments are reset simultaneously only after reaching  
the maximum episode length. In contrast, the *Partial Reset*  
mode resets each sub-environment immediately upon ter-  
mination, independent of others. Figure 7a and Figure 7b  
illustrate the differences between these two modes.

**Critic Design** Adapting PPO to VLA models places  
stringent requirements on critic design. Maintaining a sepa-  
rate critic network is computationally prohibitive and com-  
plicates GPU resource management due to the large scale of  
VLA models. We therefore adopt a parameter-sharing strat-  
egy between the actor and critic. Following RL4VLA [22],  
we attach a lightweight value head to the language model  
component of the VLA for efficient state value estimation.

**Value for Action Chunks** As discussed in Section 4.3.1,  
the two advantage estimation formulations naturally induce  
two corresponding value estimation strategies for an ac-  
tion chunk  $c_t = (a_{t,1}, a_{t,2}, \dots, a_{t,C})$ . In the *chunk-level*

<sup>1</sup>Here,  $t$  denotes the index of the chunk, and  $C$  denotes the chunk size.

374 formulation, the critic estimates a single scalar value for  
375 the entire chunk, treating it as a macro-action, i.e.,  $V : \mathcal{S} \rightarrow \mathbb{R}$ . In contrast, the *action-level* formulation pro-  
376 duces a  $C$ -dimensional value vector, providing an individ-  
377 ual value estimate for each atomic action  $a_{t,i}$  in the chunk,  
378 i.e.,  $V : \mathcal{S} \rightarrow \mathbb{R}^C$ . Empirically, we find that the action-level  
379 formulation consistently leads to better performance.  
380

### 381 4.3.3. Design Choices for GRPO

382 Transferring GRPO to embodied tasks introduces several  
383 non-trivial challenges. To address these, we adopt the fol-  
384 lowing design choices.

385 **Valid Action Mask.** While rollouts are executed for a  
386 fixed episode duration, tasks often terminate early. This  
387 presents two potential strategies for policy optimization: (1)  
388 using the full trajectory regardless of task completion tim-  
389 ing, or (2) considering only timesteps prior to task comple-  
390 tion. Our framework supports both strategies, referring to  
391 the latter as the *Valid Action Mask* setting (Figure 7c). Em-  
392 pirically, we find that masking out post-completion steps  
393 generally improves policy performance in the GRPO set-  
394 ting.

395 **Loss Normalization by Trajectory Length.** To en-  
396 sure that trajectories of different lengths contribute compa-  
397 rably to optimization, we normalize the policy loss by the  
398 number of valid timesteps in the *Valid Action Mask* setting.  
399 Specifically, for a trajectory  $\tau_i$  with  $T_i^{\text{succ}}$  valid timesteps,  
400 the contribution of each timestep to the objective is scaled  
401 by  $1/T_i^{\text{succ}}$ . This prevents longer trajectories from domi-  
402 nating the gradient and promotes balanced learning across  
403 successful and partially completed trajectories.

404 **Success Rate Filter.** Inspired by the dynamic sampling  
405 strategy in DAPO [42], we introduce a success-rate filter for  
406 GRPO. This filter discards trajectory groups where all tra-  
407 jectories either succeed or fail, as computing non-zero ad-  
408 vantages requires a mix of successful and failed outcomes.  
409 In practice, this mechanism accelerates convergence and  
410 consistently improves policy performance.

## 411 5. Experiment Results

412 In this section, we systematically investigate three key ques-  
413 tions concerning the effectiveness of the proposed RLinf-  
414 VLA framework:

415 (1) **Is RLinf-VLA high-performance?** We evaluate  
416 RLinf-VLA on three representative testbeds, LIBERO,  
417 ManiSkill and RoboTwin. The results demonstrate that  
418 RLinf-VLA achieves approximately 20–85% improvement,  
419 highlighting its strong capability to support large-scale  
420 multi-task learning.

421 (2) **Is RLinf-VLA high-efficiency?** We benchmark  
422 the framework across both GPU-parallelized and CPU-  
423 parallelized simulators, and observe that the optimal con-  
424 figuration improves training throughput, with speedups of

Table 1. Evaluation results across three simulation benchmarks. Values denote success rates (%).

(a) Evaluation on ManiSkill.							
Method	In-Distribution			OOD Avg.			
OpenVLA (Base)	53.91			39.10			
OpenVLA (RLinf-GRPO)	84.38			75.15			
OpenVLA (RLinf-PPO)	96.09			<b>81.93</b>			
OpenVLA-OFT (Base)	28.13			18.29			
OpenVLA-OFT (RLinf-GRPO)	94.14			60.64			
OpenVLA-OFT (RLinf-PPO)	<b>97.66</b>			77.05			
(b) Evaluation on LIBERO.							
OpenVLA-OFT	Object	Spatial	Goal	Long	90	Avg.	
Base	50.20	51.61	49.40	11.90	42.67	42.09	
RLinf-GRPO	99.67	98.93	98.32	93.55	98.12	98.11	
(c) Evaluation on RoboTwin.							
OpenVLA-OFT	Cup	Hammer	Bottles	Can	Pot	Hand	Avg.
Base	75.78	10.15	20.31	9.37	3.13	28.13	24.48
RLinf-GRPO	94.53	96.09	92.96	83.59	70.31	70.31	84.63

up to  $1.88\times$ . This finding underscores the necessity of  
supporting diverse allocation modes. Notably, RLinf-VLA  
achieves up to  $2.27\times$  speedup over existing frameworks.

(3) **What are the actionable practices for applying PPO and GRPO to VLA training?** Through extensive ablation studies, we identify the key factors that govern training performance, offering practical guidelines for effectively deploying RL in VLA settings.

### 5.1. Policy Performance

In this section, we evaluate the effectiveness of RLinf-VLA policies across multiple benchmarks, focusing on task success rates and generalization performance.

**Experiment Setup.** We evaluate RLinf-VLA on three benchmarks: ManiSkill, LIBERO, and RoboTwin, with full settings in Appendix III.A.

For ManiSkill, we train on 25 pick-and-place tasks with object and receptacle variations. We follow the out-of-distribution evaluation protocol of RL4VLA [22], measuring generalization in *vision*, *language*, and *action*. Each sub-setting is evaluated with 256 random episodes. For LIBERO, we use the public task groups [21], including LIBERO-Spatial, Object, Goal, Long, and 90, for training and evaluation. Instead of training within a single group, we train one unified model on the combined 130 tasks over 5 groups, testing large-scale multitask learning. Performance is reported across groups, averaged over 50 episodes per task with three random seeds. For RoboTwin [8], we select six tasks: “beat block hammer”, “move can pot”, “place empty cup”, “pick dual bottles”, “lift pot”, and “handover

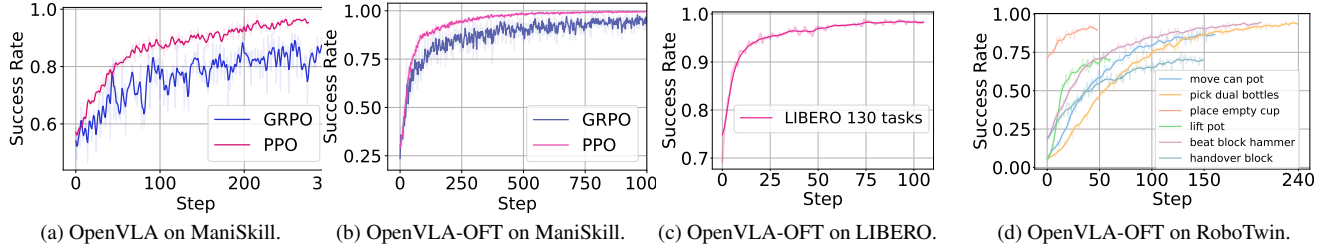


Figure 8. Training curves on different benchmarks. The x-axis shows the number of training epochs, and the y-axis indicates the corresponding success rate. Light-colored lines represent the raw data, while the dark-colored curves are smoothed using a Gaussian filter with  $\sigma = 1$ .

454 block”. Training uses 1,000 fixed randomized scene seeds  
 455 per task, and evaluation samples 128 unseen seeds to assess  
 456 out-of-distribution generalization.

457 **Training Results.** Figure 8 presents the training  
 458 curves for each setup. Figure 8a and Figure 8b report the  
 459 training performance on ManiSkill, utilizing OpenVLA and  
 460 OpenVLA-OFT models trained via PPO and GRPO. Across  
 461 all settings, RL delivers substantial performance gains, im-  
 462 proving success rates by 45%–70% over the baseline. No-  
 463 tably, PPO consistently outperforms GRPO and exhibits  
 464 greater stability for both OpenVLA and OpenVLA-OFT in  
 465 the ManiSkill setting.

466 Figure 8c illustrates the training curve of OpenVLA-  
 467 OFT with the GRPO algorithm on 130 tasks of LIBERO.  
 468 Results demonstrate that the success rate improves substan-  
 469 tially from approximately 73% to 98%, yielding an overall  
 470 performance gain of about 30%. These findings underscore  
 471 the effectiveness of the GRPO design in RLinf-VLA for en-  
 472 hancing multi-task training.

473 Figure 8d depicts our training curves for OpenVLA-OFT  
 474 on RoboTwin. The training process remains stable with  
 475 minimal fluctuations. Moreover, the model converges to  
 476 high performance even on tasks with extremely low initial  
 477 success rates (as low as 3%).

478 **Evaluation Results.** Table 1 summarizes RLinf-VLA  
 479 results across three benchmarks, with ManiSkill detailed  
 480 in Table 1a. “OpenVLA (Base)” and “OpenVLA-OFT  
 481 (Base)” denote the RL base models for OpenVLA and  
 482 OpenVLA-OFT.

483 Direct OOD comparison between OpenVLA and  
 484 OpenVLA-OFT is not strictly fair because their base per-  
 485 formance differs. OpenVLA improves from 39.10% to  
 486 81.93% after RL, while OpenVLA-OFT rises from 18.29%  
 487 to 77.05%. Although starting weaker, OpenVLA-OFT  
 488 shows slightly larger relative gains, indicating that base  
 489 model choice strongly affects final generalization.

490 Table 1b reports OpenVLA-OFT trained on 130 tasks  
 491 with RLinf-VLA over five LIBERO groups. LIBERO-  
 492 Object reaches 93%–99% success, with an overall average  
 493 of 98.11% and a mean improvement of 56.02%. Notably,  
 494 we use a single unified model across all tasks, unlike stan-

dard RL-for-VLA methods that train separate models per  
 group, demonstrating RLinf-VLA’s support for large-scale  
 multi-task RL.

OOD results on RoboTwin are shown in Table 1c. RLinf-  
 VLA achieves an average success rate of 84.63%, over 60%  
 improvement, showing strong generalization on complex  
 tasks such as long-horizon bimanual manipulation.

## 5.2. System Efficiency

In this section, we evaluate the efficiency of our system by  
 comparing different GPU allocation strategies. Our results  
 show that the optimal strategy varies across simulators and  
 models, highlighting the importance of flexible GPU allo-  
 cation, which is a core feature of RLinf-VLA.

**Experiment Setup.** We summarize key settings here,  
 with full details in Appendix III.B. We evaluate represen-  
 tative tasks from three benchmarks. For ManiSkill, we use  
 “PutCarrotOnPlateInScene-v2” with 256 parallel environ-  
 ments. For LIBERO, we adopt the LIBERO-Long suite, and  
 for RoboTwin, the “place empty cup” task. For LIBERO  
 and RoboTwin, parallel environments scale linearly with  
 nodes (64, 128, and 256 for 1, 2, and 4 nodes).

We consider three GPU allocation strategies (Sec-  
 tion 4.1). *Collocated* shares GPUs across components,  
*Disaggregated* separates training from rollout, and *Hy-  
 brid* pipelines simulation by splitting instances into stages.  
 For comparisons, ManiSkill lacks multi-GPU baselines,  
 so we compare *Collocated* and *Hybrid* (pipelined) modes  
 against a naive *Disaggregated* baseline. For LIBERO and  
 RoboTwin, we benchmark against SimpleVLA-RL [20]  
 (collocated only) and evaluate our *Collocated* and *Hybrid*  
 modes; the disaggregated mode is omitted since hybrid  
 consistently performs better when GPUs are underutilized.

**Results.** Figure 9 reports the end-to-end throughput of  
 RLinf-VLA and baselines for VLA+RL training under dif-  
 ferent cluster sizes and placement strategies.

From Figure 9a, for OpenVLA in ManiSkill, RLinf-  
 VLA achieves up to a  $1.88\times$  speedup over the disaggre-  
 gated baseline using the Hybrid (pipe=2) mode on 8  
 GPUs. Although scaling to more GPUs introduces over-  
 heads from model loading, offloading, and state switch-

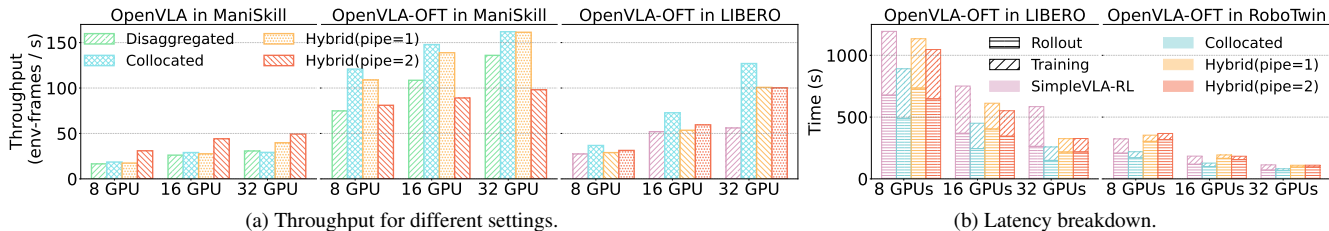


Figure 9. Evaluation of system efficiency on ManiSkill, LIBERO, and RoboTwin, comparing OpenVLA and OpenVLA-OFT. Throughput (total environment frames per second) improves with increasing pipeline stages (“pipe”).

535 ing, fine-grained pipelining reduces GPU idle time and  
 536 maintains a  $1.61\times-1.69\times$  advantage over the disaggregated  
 537 mode. Deeper pipelines (pipe=2 vs. pipe=1) further im-  
 538 prove performance by reducing rollout-stage bubbles.

539 For OpenVLA in ManiSkill and OpenVLA-OFT in  
 540 ManiSkill, as GPU count increases, the hybrid (pipe=1)  
 541 mode outperforms the collocated mode. Since ManiSkill  
 542 is GPU-parallelized, rollout throughput scales with parallel  
 543 environments, making simulator resource allocation critical.  
 544 To improve utilization, we enable offloading in the collocated  
 545 mode (Section 4.1.1), but communication over-  
 546 head grows with cluster size. In contrast, hybrid mode splits  
 547 GPUs evenly between components, reducing interference.

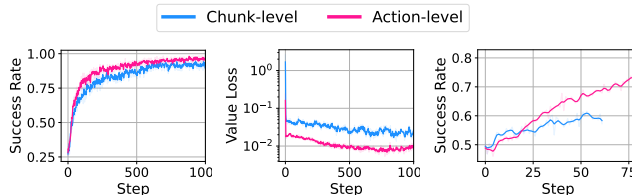
548 For OpenVLA-OFT in LIBERO and ManiSkill, the behav-  
 549 ior differs: collocated and hybrid (pipe=1) modes  
 550 outperform disaggregated and hybrid (pipe=2). First,  
 551 OpenVLA-OFT generates action chunks while the simula-  
 552 tor executes them sequentially, shifting the generation-to-  
 553 execution ratio from 1:1 to about 1:15. Second, LIBERO’s  
 554 CPU-parallelized simulator becomes the main bottleneck.  
 555 This imbalance diminishes the benefit of pipelining, effec-  
 556 tively reverting execution to a near-sequential process.

557 From Figure 9b, where “Rollout” denotes the total time  
 558 for multi-step interactions between the generator and simu-  
 559 lator, for OpenVLA-OFT in LIBERO and RoboTwin,  
 560 RLinf-VLA still achieves a  $1.34\times-2.27\times$  speedup over  
 561 SimpleVLA-RL under collocated settings. The gains come  
 562 from both rollout and training. During rollout, RLinf-VLA  
 563 uses vectorized environments that are more efficient than  
 564 SimpleVLA-RL’s multi-process workers. During training,  
 565 system-level optimizations such as adaptive communication  
 566 and the removal of redundant log-probability recomputation  
 567 further reduce latency, with benefits increasing at scale.

568 **Summary.** Different simulators demand different allo-  
 569 cation strategies. RLinf-VLA’s flexible multi-mode support  
 570 enables consistently high efficiency across diverse execu-  
 571 tion regimes.

### 572 5.3. Ablation Study

573 In this section, we present ablation studies under differ-  
 574 ent setups and summarize key insights for PPO and GRPO  
 575 training.

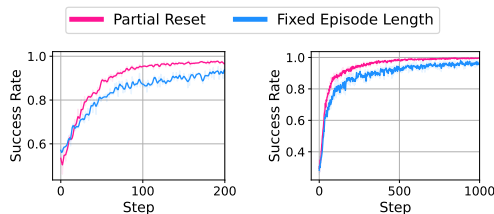


(a) succ. in ManiSkill (b) loss in ManiSkill (c) succ. in LIBERO

Figure 10. For OpenVLA-OFT, action-level value estimation consistently outperforms chunk-level estimation across different tasks.

#### (1) Tips for PPO.

576 (a) *Action-level value estimation outperforms chunk-*  
 577 *level estimation for PPO with action chunks.* Figure 10  
 578 shows the PPO training curves comparing action-level and  
 579 chunk-level value estimation as described in Section 4.3.2  
 580 on the ManiSkill “PutOnPlateInScene25Main” task using  
 581 the OpenVLA-OFT model. Action-level estimation consis-  
 582 tently yields higher success rates and lower value loss,  
 583 demonstrating more effective learning and faster policy im-  
 584 provement. The performance divergence between different  
 585 levels of value estimation remains consistent across diverse  
 586 tasks, as corroborated by similar results in the LIBERO-  
 587 Goal benchmark (Figure 10c).  
 588



(a) succ. in OpenVLA (b) succ. in OpenVLA-OFT

Figure 11. Partial reset consistently outperforms Fixed Episode Length mode in PPO.

589 (b) *Partial reset substantially improves sample efficiency.*  
 590 Figure 11 presents the PPO training curves for the *Partial*  
 591 *Reset* and *Fixed Episode Length* rollout modes discussed in  
 592 Section 4.3.2. Since the optimization objective is to achieve  
 593 success at least once per rollout (“success\_once”), *Partial*  
 594 *Reset* leads to a significantly higher success rate. For a given  
 595

595 number of training steps, the success rate under *Partial Re-*  
 596 *set* consistently exceeds that of the *Fixed Episode Length*  
 597 mode. This trend is observed regardless of the model type  
 598 (OpenVLA or OpenVLA-OFT).

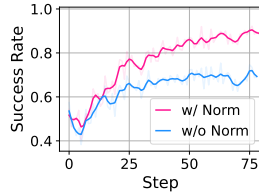


Figure 12. Ablation study on trajectory length normalization. Success rate in LIBERO-Goal with OpenVLA-OFT.

## 599 (2) Tips for GRPO.

600 (a) *Trajectory length normalization in GRPO.* As dis-  
 601 cussed in Section 4.3.3, normalizing the loss by trajectory  
 602 length aims to reduce bias when episodes vary in length.  
 603 Figure 12 shows that incorporating trajectory length nor-  
 604 malization (“w/ Norm”) can lead to substantially higher  
 605 performance compared with the unnormalized setting (“w/o  
 606 Norm”).

607 (b) *Valid action mask in GRPO.* We also study the ef-  
 608 fect of valid action masking, introduced in Section 4.3.3.  
 609 Since the optimization objective is defined with respect to  
 610 “success\_once”, excluding actions after reaching the suc-  
 611 cess state improves sample efficiency and avoids redundant  
 612 updates. Moreover, applying the valid action mask natu-  
 613 rally results in shorter trajectories, which further benefits  
 614 trajectory length normalization. Figure 13a shows results  
 615 on LIBERO-Goal. Comparing the “w/o Mask, w/o Norm”  
 616 and “w/ Mask, w/o Norm” curves, the setting with a valid  
 617 action mask achieves consistently better performance. The  
 618 “w/ Mask, w/ Norm” curve demonstrates that combining  
 619 the mask with trajectory length normalization provides an  
 620 additional improvement. However, the effect of these tech-  
 621 niques is task-dependent. In the ManiSkill setting, we do  
 622 not observe clear benefits from either valid action masking  
 623 or trajectory length normalization.

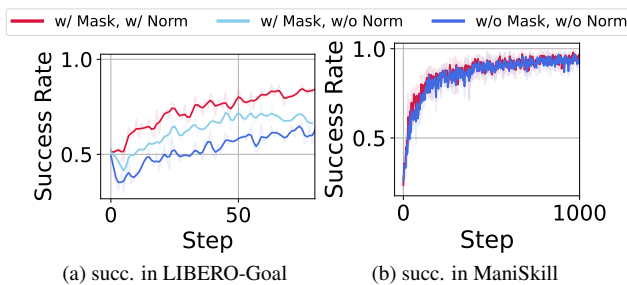


Figure 13. Ablation studies on valid action mask in GRPO with OpenVLA-OFT.

624 (c) *Success rate filtering can improve training stability*  
 625 *in some settings of GRPO.* As discussed in Section 4.3.3,

we implement a success rate filter for GRPO that discards  
 groups in which all trajectories have identical cumulative  
 rewards. This mechanism can improve the stability of  
 GRPO training. For instance, in the OpenVLA ManiSkill  
 setting, training without the filter (“w/o Filter”) exhibits a  
 clear collapse around step 400, whereas enabling the fil-  
 ter (“w/ Filter”) largely alleviates this issue. However, the  
 benefit of the filter is not universal: in the OpenVLA-OFT  
 ManiSkill setting and the OpenVLA-OFT LIBERO-Goal  
 setting, its effectiveness is much less pronounced.

## 6. Conclusion

In this work, we introduced RLinf-VLA, a unified and ef-  
 ficient framework for reinforcement learning-based train-  
 ing of Vision-Language-Action models. RLinf-VLA inte-  
 grates multiple simulators, algorithms, and VLA architec-  
 tures, while providing flexible execution modes and system-  
 level optimizations that significantly improve training ef-  
 ficiency. Extensive experiments demonstrate that models  
 trained with RLinf-VLA achieve approximately 20–85%  
 improvement on a wide range of simulated tasks. More-  
 over, our study distills actionable practices for both PPO  
 and GRPO, guiding future research in RL-based VLA train-  
 ing. By open-sourcing RLinf-VLA with ongoing mainte-  
 nance, we provide the community with a foundation to ac-  
 celerate, standardize, and scale research in embodied intel-  
 ligence.

## References

- [1] Hongzhe Bi, Hengkai Tan, Shenghao Xie, Zeyuan Wang, Shuhe Huang, Haitian Liu, Ruowen Zhao, Yao Feng, Chendong Xiang, Yinze Rong, Hongyan Zhao, Hanyu Liu, Zhizhong Su, Lei Ma, Hang Su, and Jun Zhu. Motus: A unified latent action world model. *arXiv preprint arXiv:2512.13030*, 2025. 2
- [2] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ : A vision-language-action flow model for general robot control, 2024. 2
- [3] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xu Huang, Shu Jiang, et al. Agibot world colosseum: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025. 3
- [4] Qingwen Bu, Yanting Yang, Jisong Cai, Shenyuan Gao, Guanghui Ren, Maoqing Yao, Ping Luo, and Hongyang Li. Univla: Learning to act anywhere with task-centric latent actions. *arXiv preprint arXiv:2505.06111*, 2025. 2
- [5] Jun Cen, Siteng Huang, Yuqian Yuan, Kehan Li, Hangjie Yuan, Chaohui Yu, Yuming Jiang, Jiayan Guo, Xin Li, Hao

- 678 Luo, Fan Wang, Fan Wang, and Deli Zhao. Rynnvla-002:  
679 A unified vision-language-action and world model. *arXiv*  
680 *preprint arXiv:2511.17502*, 2025.
- 681 [6] Chilam Cheang, Sijin Chen, Zhongren Cui, Yingdong Hu,  
682 Liqun Huang, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu,  
683 Xiao Ma, et al. Gr-3 technical report. *arXiv preprint*  
684 *arXiv:2507.15493*, 2025. 2
- 685 [7] Kang Chen, Zhihao Liu, Tonghe Zhang, Zhen Guo, Si  
686 Xu, Hao Lin, Hongzhi Zang, Quanlu Zhang, Zhaofei Yu,  
687 Guoliang Fan, et al.  $\pi_{r1}$ : Online rl fine-tuning for  
688 flow-based vision-language-action models. *arXiv preprint*  
689 *arXiv:2510.25889*, 2025. 2
- 690 [8] Tianxing Chen, Zanxin Chen, Baijun Chen, Zijian Cai, Yibin  
691 Liu, Zixuan Li, Qiwei Liang, Xianliang Lin, Yiheng Ge,  
692 Zhenyu Gu, Weiliang Deng, Yubin Guo, Tian Nian, Xuan-  
693 bing Xie, Qiangyu Chen, Kailun Su, Tianling Xu, Guodong  
694 Liu, Mengkang Hu, Huan ang Gao, Kaixuan Wang, Zhix-  
695 uan Liang, Yusen Qin, Xiaokang Yang, Ping Luo, and Yao  
696 Mu. Robotwin 2.0: A scalable data generator and bench-  
697 mark with strong domain randomization for robust bimanual  
698 robotic manipulation, 2025. 2, 3, 5, 6, 15
- 699 [9] Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha  
700 Majumdar, Jiankai Sun, Weiyu Liu, Yuke Zhu, Shuran Song,  
701 Ashish Kapoor, Karol Hausman, et al. Foundation models  
702 in robotics: Applications, challenges, and the future. *The In-*  
703 *ternational Journal of Robotics Research (IJRR)*, 44(5):701–  
704 739, 2025. 2
- 705 [10] Haoran Geng, Feishi Wang, Songlin Wei, Yuyang Li,  
706 Bangjun Wang, Boshi An, Charlie Tianyue Cheng, Haozhe  
707 Lou, Peihao Li, Yen-Jen Wang, et al. Roboverse: To-  
708 wards a unified platform, dataset and benchmark for scal-  
709 able and generalizable robot learning. *arXiv preprint*  
710 *arXiv:2504.18904*, 2025. 3
- 711 [11] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,  
712 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi  
713 Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning  
714 capability in llms via reinforcement learning. *arXiv preprint*  
715 *arXiv:2501.12948*, 2025. 5
- 716 [12] Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-  
717 Jen Wang, Yucheng Hu, and Jianyu Chen. Improving vision-  
718 language-action model with online reinforcement learning.  
719 *IEEE International Conference on Robotics and Automation*  
720 *(ICRA)*, 2025. 2
- 721 [13] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-  
722 Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora:  
723 Low-rank adaptation of large language models. In *Inter-*  
724 *national Conference on Learning Representations (ICLR)*,  
725 2022. 5, 15, 18
- 726 [14] Physical Intelligence, Ali Amin, Raichelle Aniceto, Ash-  
727 win Balakrishna, Kevin Black, Ken Conley, Grace Con-  
728 nors, James Darpinian, Karan Dhabalia, Jared DiCarlo, et al.  
729  $\pi_{0.6}^*$ : a vla that learns from experience. *arXiv preprint*  
730 *arXiv:2511.14759*, 2025. 2
- 731 [15] Physical Intelligence, Kevin Black, Noah Brown, James  
732 Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail,  
733 Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y.  
734 Galliker, Dibya Ghosh, et al.  $\pi_{0.5}$ : a vision-language-  
action model with open-world generalization. *arXiv preprint*  
*arXiv:2504.16054*, 2025. 2
- [16] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Bal-  
akrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush  
Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang  
Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild  
robot manipulation dataset. In *Robotics: Science and Sys-*  
*tems (RSS)*, 2024. 2
- [17] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao,  
Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Fos-  
ter, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kol-  
lar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey  
Levine, Percy Liang, and Chelsea Finn. Openvla: An open-  
source vision-language-action model. *Conference on Robot*  
*Learning (CoRL)*, 2024. 2, 3, 5, 15
- [18] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning  
vision-language-action models: Optimizing speed and suc-  
cess. *arXiv preprint arXiv:2502.19645*, 2025. 2, 3, 5, 15
- [19] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen,  
Sanjana Srivastava, Roberto Martín-Martín, Chen Wang,  
Gabrael Levine, Wensi Ai, Benjamin Jose Martinez, et al.  
Behavior-1k: A human-centered, embodied ai benchmark  
with 1, 000 everyday activities and realistic simulation.  
*CoRR*, abs/2403.09227, 2024. 3
- [20] 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 train-  
ing via reinforcement learning. 2026. 2, 7, 15, 16
- [21] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu,  
Yuke Zhu, and Peter Stone. Libero: Benchmarking knowl-  
edge transfer for lifelong robot learning. *Advances in Neural*  
*Information Processing Systems*, 36:44776–44791, 2023. 2,  
3, 5, 6, 15
- [22] 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. In *The Thirty-*  
*ninth Annual Conference on Neural Information Processing*  
*Systems*, 2025. 2, 5, 6, 15
- [23] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan,  
Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu.  
Rdt-1b: A diffusion foundation model for bimanual manip-  
ulation. *International Conference on Learning Representa-*  
*tions (ICLR)*, 2025. 2
- [24] Guanxing Lu, Wenkai Guo, Chubin Zhang, Yuheng Zhou,  
Haonan Jiang, Zifeng Gao, Yansong Tang, and Ziwei Wang.  
Vla-rl: Towards masterful and general robotic manipula-  
tion with scalable reinforcement learning. *arXiv preprint*  
*arXiv:2505.18719*, 2025. 2
- [25] Yuen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and  
Irwin King. A survey on vision-language-action models for  
embodied ai. *arXiv preprint arXiv:2405.14093*, 2024. 2
- [26] Reginald McLean, Evangelos Chatzaroulas, Luc Mc-  
Cutcheon, Frank Röder, Tianhe Yu, Zhanpeng He, K.R.  
Zentner, Ryan Julian, J K Terry, Isaac Woungang, Nariman  
Farsad, and Pablo Samuel Castro. *Meta-World+*: An Im-  
proved, Standardized, RL Benchmark. *Advances in Neu-*  
*ral Information Processing Systems (NeurIPS) Datasets and*  
*Benchmarks Track*, 2025. 3

- 793 [27] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wol- 850  
794 fram Burgard. Calvin: A benchmark for language- 851  
795 conditioned policy learning for long-horizon robot manipu- 852  
796 lation tasks. *IEEE Robotics and Automation Letters (RA-L)*, 853  
797 7(3):7327–7334, 2022.
- 798 [28] Mayank Mittal, Pascal Roth, James Tigue, Antoine Richard, 854  
799 Octi Zhang, Peter Du, Antonio Serrano-Muñoz, Xinjie Yao, 855  
800 René Zurbrügg, et al. Isaac lab: A gpu-accelerated simu- 856  
801 lation framework for multi-modal robot learning. *arXiv* 857  
802 *preprint arXiv:2511.04831*, 2025. 858
- 803 [29] Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, 859  
804 Adeet Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, 860  
805 and Yuke Zhu. *RoboCasa: Large-Scale Simulation of Ev-* 861  
806 *eryday Tasks for Generalist Robots*. Robotics: Science and 862  
807 Systems (RSS), 2024. 3 863
- 808 [30] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Ab- 864  
809 hishek Gupta, Abhishek Padalkar, Abraham Lee, et al. Open 865  
810 x-embodiment: Robotic learning datasets and rt-x models. 866  
811 *IEEE International Conference on Robotics and Automation* 867  
812 *(ICRA)*, pages 6892–6903, 2024. 2 868
- 813 [31] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Er- 869  
814 mon, Christopher D. Manning, and Chelsea Finn. Direct 870  
815 preference optimization: Your language model is secretly a 871  
816 reward model. *Advances in Neural Information Processing* 872  
817 *Systems (NeurIPS)*, 2023. 2 873
- 818 [32] John Schulman, Philipp Moritz, Sergey Levine, Michael I. 874  
819 Jordan, and P. Abbeel. High-dimensional continuous 875  
820 control using generalized advantage estimation. *CoRR*, 876  
821 abs/1506.02438, 2015. 877
- 822 [33] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Rad- 878  
823 ford, and Oleg Klimov. Proximal policy optimization algo- 879  
824 rithms. *arXiv preprint arXiv:1707.06347*, 2017. 2, 5, 12 880
- 825 [34] Dhruv Shah, Błażej Osiniński, Sergey Levine, et al. Lm- 881  
826 nav: Robotic navigation with large pre-trained models of 882  
827 language, vision, and action. *Conference on Robot Learn-* 883  
828 *ing (CoRL)*, pages 492–504, 2022. 2 884
- 829 [35] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao 885  
830 Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, 886  
831 Yang Wu, et al. Deepseekmath: Pushing the limits of math- 887  
832 ematical reasoning in open language models. *arXiv preprint* 888  
833 *arXiv:2402.03300*, 2024. 2, 3, 5, 13 889
- 834 [36] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, 890  
835 Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and 891  
836 Chuan Wu. Hybridflow: A flexible and efficient rlhf frame- 892  
837 work. In *Proceedings of the Twentieth European Confer-* 893  
838 *ence on Computer Systems*, page 1279–1297, New York, NY, 894  
839 USA, 2025. Association for Computing Machinery. 2 895
- 840 [37] Richard S Sutton and Andrew G Barto. *Reinforcement* 896  
841 *Learning: An Introduction*. MIT Press, 1998. 2 897
- 842 [38] Shuhan Tan, Kairan Dou, Yue Zhao, and Philipp 898  
843 Krähenbühl. Interactive post-training for vision-language- 899  
844 action models. *arXiv preprint arXiv:2505.17016*, 2025. 2 900
- 845 [39] Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander 901  
846 Hinrichsen, Xiaodi Yuan, Chen Bao, Xinsong Lin, Yulin Liu, 902  
847 Tse kai Chan, Yuan Gao, Xuanlin Li, Tongzhou Mu, Nan 903  
848 Xiao, Arnav Gurha, Viswesh Nagaswamy Rajesh, Yong Woo 904  
849 Choi, Yen-Ru Chen, Zhiao Huang, Roberto Calandra, Rui 905  
Chen, Shan Luo, and Hao Su. *ManiSkill3: GPU Parallelized* 906  
*Robotics Simulation and Rendering for Generalizable Em-* 907  
*bodied AI*. Robotics: Science and Systems, 2025. 2, 3, 5, 13, 908  
15 909
- [40] Octo Model Team, Dibya Ghosh, Homer Walke, Karl 910  
Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey 911  
Hejna, Tobias Kreiman, Charles Xu, et al. *Octo: An Open-* 912  
*Source Generalist Robot Policy*. Robotics: Science and Sys- 913  
tems, 2024. 2 914
- [41] Junjie Wen, Yichen Zhu, Jinming Li, Zhibin Tang, Chaomin 915  
Shen, and Feifei Feng. Dexvla: Vision-language model with 916  
plug-in diffusion expert for general robot control. *arXiv* 917  
*preprint arXiv:2502.05855*, 2025. 2 918
- [42] Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xi- 919  
aochen Zuo, YuYue, Weinan Dai, Tiantian Fan, Gaohong 920  
Liu, Juncai Liu, LingJun Liu, Xin Liu, Haibin Lin, Zhiqi 921  
Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, 922  
Mofan Zhang, Ru Zhang, Wang Zhang, Hang Zhu, Jinhua 923  
Zhu, Jiaye Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, 924  
Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei- 925  
Ying Ma, Ya-Qin Zhang, Lin Yan, Yonghui Wu, and Mingx- 926  
uan Wang. DAPO: An open-source LLM reinforcement 927  
learning system at scale. In *The Thirty-ninth Annual Con-* 928  
*ference on Neural Information Processing Systems*, 2025. 6 929
- [43] Zhecheng Yuan, Tianming Wei, Shuiqi Cheng, Gu Zhang, 930  
Yuanpei Chen, and Huazhe Xu. Learning to manipulate any- 931  
where: A visual generalizable framework for reinforcement 932  
learning. *arXiv preprint arXiv:2407.15815*, 2024. 2 933
- [44] Zijian Zhang, Kaiyuan Zheng, Zhaorun Chen, Joel Jang, Yi 934  
Li, Siwei Han, Chaoqi Wang, Mingyu Ding, Dieter Fox, and 935  
Huaxiu Yao. Grape: Generalizing robot policy via prefer- 936  
ence alignment. *arXiv preprint arXiv:2411.19309*, 2024. 2 937
- [45] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea 938  
Finn. Learning fine-grained bimanual manipulation with 939  
low-cost hardware. In *Robotics: Science and Systems (RSS)*, 940  
pages 1–16. Robotics: Science and Systems, 2023. 3 941
- [46] Jinliang Zheng, Jianxiong Li, Zhihao Wang, Dongxiu Liu, 942  
Xirui Kang, Yuchun Feng, Yanan Zheng, Jiayin Zou, Yilun 943  
Chen, Jia Zeng, et al. X-vla: Soft-prompted transformer 944  
as scalable cross-embodiment vision-language-action model. 945  
*arXiv preprint arXiv:2510.10274*, 2025. 2 946

# Appendix for RLinf-VLA

891

## 7. Theoretical Background

### 7.1. Definitions of RL

**Definition 1** (Return). The return  $R_t$  is the  $\gamma$ -discounted cumulative reward starting from timestep  $t$ :

$$R_t = \sum_{k=0}^{\infty} \gamma^k r(s_{t+k}, a_{t+k}). \quad (2)$$

**Definition 2** (Value Function). The value function  $V_\pi(s)$  is the expected return when starting from state  $s$  and following policy  $\pi$ :

$$V_\pi(s) = \mathbb{E}_\pi[R_t \mid s_t = s], \quad (3)$$

where  $t$  denotes an arbitrary timestep at which the agent is in state  $s$ .

**Definition 3** (Action-Value Function). The action-value function  $Q_\pi(s, a)$  is the expected return when executing action  $a$  in state  $s$  and thereafter following policy  $\pi$ :

$$Q_\pi(s, a) = \mathbb{E}_\pi[R_t \mid s_t = s, a_t = a], \quad (4)$$

where  $t$  denotes an arbitrary timestep at which the agent is in state  $s$  and takes action  $a$ .

**Definition 4** (Advantage Function). The advantage function  $A_\pi(s, a)$  quantifies how much better taking action  $a$  in state  $s$  is compared to the expected value of the state:

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(s). \quad (5)$$

### 7.2. Introduction of PPO

PPO [33] is one of the most widely adopted reinforcement learning algorithms in robotics. PPO enhances training stability by constraining policy updates within a trust region, thereby preventing overly large changes that could destabilize learning. This is achieved through a clipped surrogate objective, which balances exploration and exploitation while maintaining sample efficiency. Due to its robustness and simplicity, PPO has become a standard baseline in robotics RL and serves as the foundation of our framework.

In PPO, the advantage function is commonly estimated using Generalized Advantage Estimation (GAE) [32]:

$$\hat{A}_t = \sum_{k=0}^{T-t-1} (\gamma\lambda)^k (r_{t+k} + \gamma V(s_{t+k+1}) - V(s_{t+k})), \quad (6)$$

where  $r_t$  denotes the reward at timestep  $t$ ,  $V(s)$  is the value function,  $\gamma$  is the discount factor,  $\lambda$  controls the bias–variance trade-off, and  $T$  is the episode horizon.

The PPO optimization objective is defined as:

$$J^{\text{PPO}}(\theta) = \mathbb{E}_t \left[ \min(\rho_t(\theta)\hat{A}_t, \text{clip}(\rho_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right], \quad (7)$$

where

$$\rho_t(\theta) = \frac{\pi_\theta(a_t \mid o_t)}{\pi_{\theta_{\text{old}}}(a_t \mid o_t)}, \quad (8)$$

$\pi_{\theta_{\text{old}}}$  denotes the rollout policy, and  $\epsilon$  is the clipping parameter.<sup>2</sup>

<sup>2</sup>The definition of action  $a_t$  differs across LLM and robotics RL settings. Here we treat it as a generalized action. The same applies to GRPO.

### 7.3. Introduction of GRPO

GRPO [35] is a recent variant of policy optimization designed to simplify reinforcement learning pipelines. Unlike PPO, which requires training both a policy model and a value model, GRPO eliminates the need for an explicit value function by leveraging group-based relative comparisons of trajectories. This design reduces the overall model complexity and avoids potential inaccuracies from value estimation. As a result, GRPO provides a lightweight yet effective alternative to PPO, making it especially attractive for large-scale VLA training where efficiency and simplicity are important. In GRPO, the important ratio can be defined as:

$$\rho_t^{(i)}(\theta) = \frac{\pi_\theta(a_t^{(i)} | o_t^{(i)})}{\pi_{\theta_{\text{old}}}(a_t^{(i)} | o_t^{(i)})}, \quad (9)$$

where  $i$  denotes  $i$ -th trajectory. But the advantage can be defined as:

$$\hat{A}^{(i)} = \hat{A}_t^{(i)} = \frac{\mathcal{R}^{(i)} - \text{mean}(\{\mathcal{R}^{(j)}\}_{j=1}^G)}{\text{std}(\{\mathcal{R}^{(j)}\}_{j=1}^G)}, \quad \forall t \in \{0, 1, \dots, |\tau^{(i)}|\}, \quad (10)$$

where  $\mathcal{R}^{(i)}$  denotes the total reward of trajectory  $\tau^{(i)}$ . Note that  $\mathcal{R}$  differs from the discounted return  $R$ , and  $G$  represents the group size. Given an initial observation  $o_0$  from sample buffer  $D$ , the behavior model  $\pi_{\theta_{\text{old}}}$  generates  $G$  trajectories  $\{\tau_i\}_{i=1}^G$ . The GRPO optimization objective is:

$$J^{\text{GRPO}}(\theta) = \mathbb{E}_{o_0 \sim D, \{\tau^{(i)}\} \sim \pi_{\theta_{\text{old}}}} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|\tau^{(i)}|} \sum_{t=1}^{|\tau^{(i)}|} \min\left(\rho_t^{(i)}(\theta) \hat{A}^{(i)}, \text{clip}\left(\rho_t^{(i)}(\theta), 1 - \epsilon, 1 + \epsilon\right) \hat{A}^{(i)}\right) \right], \quad (11)$$

where  $|\tau^{(i)}|$  is the length of trajectory  $\tau^{(i)}$ ,  $\epsilon$  is the clip parameter.

## 8. Implementation Details of RLinf-VLA

### 8.1. The Unified Interface

The unified interface consists of two categories: *core functions* and *utility functions*.

**Core Functions** We implement standard Gym-style APIs, including “reset” and “step”. The “step” function additionally supports an “auto\_reset” option: when enabled, any sub-environment that terminates or is truncated will be automatically reset, thereby improving sample efficiency by avoiding idle environments. Following the ManiSkill [39] convention, we also support the “ignore\_terminations” option. When enabled, termination signals are ignored and only truncation signals are respected, meaning that an episode ends only when the maximum episode length is reached. This feature allows us to flexibly support different implementation variants, such as “Partial Reset” and “Valid Action Mask”.

Beyond these, we extend the interface with a “chunk\_step” function to handle action chunks. Instead of naively looping over the chunk with repeated “step” calls, this function manages episode termination more carefully. Two modes are supported: (1) reset immediately when a sub-environment finishes within the chunk, or (2) defer reset until the entire chunk has been executed. This flexibility ensures the correct handling of episode boundaries when chunked actions are used.

**Utility Functions** Utility functions offer convenient support for training and evaluation. For instance, visualization utilities enable effortless generation of videos during rollouts or evaluation. In addition, we provide specialized utilities required by specific algorithms. For example, GRPO necessitates that all environments within a group share the same initial state, which can be ensured by setting “use\_fixed\_reset\_state\_ids=True”.

### 8.2. GPU Allocation Modes

**Configuring GPU Allocation Modes** Our framework exposes a simple configuration interface that allows users to flexibly choose GPU allocation strategies without explicitly specifying the mode. Instead, users only need to assign GPU IDs for each component via the following fields:

- `cluster.component_placement.env` for the *Simulator*,
- `cluster.component_placement.rollout` for *Generation*,

- 960 • `cluster.component_placement.actor` for *Train*.
- 961 In addition, offloading can be enabled or disabled independently for each component using:
- 962 • `env.enable_offload`,
- 963 • `rollout.enable_offload`,
- 964 • `actor.enable_offload`.

965 Finally, fine-grained pipelining is controlled by the parameter `rollout.pipeline_stage_num`. A value greater than  
966 1 splits each simulator instance on a GPU into multiple pipeline stages, while a value of 1 disables pipelining.

### 967 8.3. Details for Design Choices

#### 968 8.3.1. Details for Multi-Granularity Support

969 In the main text, we refer to three levels of granularity: *token-level*, *action-level*, and *chunk-level*. Figure 14 provides an  
970 illustrative example of these levels. An *action chunk* consists of multiple atomic actions, which are the actual control signals  
971 executed by the simulator at each step. Each atomic action, in turn, is composed of multiple tokens, where each token  
972 corresponds to a single action dimension.

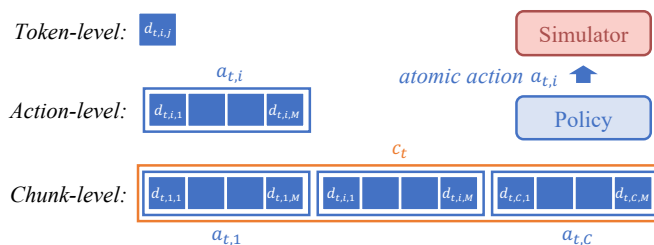


Figure 14. Illustration of different log-probability granularities.

973 The complete set of supported combinations between advantage computation granularity and log-probability granularity  
974 is summarized in Table 2.

Table 2. Supported combinations of advantage and log-probability granularities.

Advantage	Log-Probability		
	Chunk-level	Action-level	Token-level
Chunk-level	✓	✓	✓
Action-level	✗	✓	✓

#### 975 8.3.2. Details for Loss Normalization by Trajectory Length

976 As mentioned in the main text, we normalize the policy loss by the trajectory length under the *Valid Action Mask* setting.  
977 Specifically, if a trajectory  $\tau_i$  has  $T_i^{\text{succ}}$  valid timesteps, the contribution of each timestep to the objective is scaled by  $1/T_i^{\text{succ}}$ .  
978 The GRPO objective under this setting is:

$$979 J^{\text{GRPO}}(\theta) = \mathbb{E}_{o_0 \sim D, \{\tau_i\} \sim \pi_{\theta, \text{old}}} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{T_i^{\text{succ}}} \sum_{t=1}^{T_i^{\text{succ}}} \min \left( \rho_{i,t}(\theta) \hat{A}_i, \text{clip}(\rho_{i,t}(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right) \right]. \quad (12)$$

980 This prevents longer trajectories from disproportionately dominating the gradient and ensures balanced learning across tra-  
981 jectories of varying lengths.

## 982 9. Experiment Setup

### 983 9.1. Details for Performance Experiments

984 This part is the details of training and evaluation for High-Performance experiments, including metric for evaluation, base  
985 models selection and hyperparameter settings for training. Unless otherwise specified, all curves are smoothed using a  
986 Gaussian filter ( $\sigma = 1$ ), and success rate is computed under the “`success_once`” criterion, where an episode is considered  
987 successful if the success state is reached at least once.

Table 3. Hyperparameter settings for training.

Parameter	OpenVLA [17] ManiSkill [39]		OpenVLA-OFT [18] ManiSkill		OpenVLA-OFT LIBERO [21]	OpenVLA-OFT RoboTwin [8]
	PPO	GRPO	PPO	GRPO	GRPO	GRPO
# Parallel Envs	128	256	128	256	64	128
Max Episode Steps	80	80	80	80	512	200 <sup>a</sup>
Max Steps per Rollout Epoch	160	80	160	80	512	200 <sup>a</sup>
Group Size	1	8	1	8	8	8
Rollout Epoch	1	1	1	1	64	8
Partial Reset	True	False	True	False	False	False
Valid Action Mask	False	True	False	True	True	True
Action Chunk Size	1	1	8	8	8	25
Global Batch Size	640	640	640	640	16384	1024
Micro Batch Size	40	40	40	40	32	32
is_lora [13]	True	True	True	True	—	True
Learning Rate	1e-4	1e-5	1e-4	1e-4	2e-5	1e-4
adam_eps	1e-8	1e-8	1e-8	1e-8	1e-8	1e-5
Clip Ratio ( $\epsilon$ )	(0.2, 0.28)	(0.2, 0.28)	(0.2, 0.28)	(0.2, 0.28)	(0.2, 0.28)	(0.2, 0.28)
$\gamma$	0.99	—	0.99	—	—	—
GAE $\lambda$	0.95	—	0.95	—	—	—
Temperature (train)	1.0	1.0	1.0	1.0	1.6	1.6

<sup>a</sup> Note that the “pick dual bottles” and “handover block” task in RoboTwin requires 100 and 400 steps respectively.

### 9.1.1. Base Models for Training

For the OpenVLA (Base) in ManiSkill, we adopt the pre-trained checkpoint OpenVLA (Base) from RL4VLA [22]. For the OpenVLA-OFT (Base) in ManiSkill, we perform our own LoRA fine-tuning to get OpenVLA-OFT (Base) using motion planning data collected from the “PutOnPlateInScene25Main-v3” task.

For OpenVLA-OFT (Base) in LIBERO 130 tasks, we fine-tune OpenVLA-OFT with LoRA using demonstrations from all five task suites in LIBERO.

For OpenVLA-OFT (Base) in RoboTwin, we adopt the pre-trained checkpoint OpenVLA-OFT (Base) from SimpleVLA-RL [20]. Besides, we also reuse the train/eval seeds from SimpleVLA-RL [20].

### 9.1.2. Hyperparameter Settings for Training

In Table 3, we list the specific hyperparameters used for different experiments.

### 9.1.3. Metric for Evaluation

For supervised fine-tuned models, we set `do_sample=False` during evaluation. And for RL-trained models, we set `do_sample=True` during evaluation, and conduct three times evaluation for each model. And we report the mean and standard deviation of the success rates. (Now only the LIBERO setting has standard deviation.)

## 9.2. Details for Efficiency Experiments

Table 4. The setting of `enable_offload` among different GPU allocation modes.

	Simulator	Generation	Training
Disaggregated	False	False	False
Colocated	True	True	True
Hybrid	True	True	True

### 9.2.1. Hyperparameters for Efficiency Experiments

See Table 4 for how we enable the offload configurations and Table 5 for how we configure the GPU allocation.

Table 5. Hyperparameters of GPU allocation modes.

# GPU	Allocation Mode	Simulator	Generation	Training
8 GPUs	Disaggregated	0-1	2-3	4-7
	Colocated	0-7	0-7	0-7
	Hybrid	0-3	4-7	0-7
16 GPUs	Disaggregated	0-3	4-7	8-15
	Colocated	0-15	0-15	0-15
	Hybrid	0-7	8-15	0-15
32 GPUs	Disaggregated	0-7	8-15	16-31
	Colocated	0-31	0-31	0-31
	Hybrid	0-15	16-31	0-31

### 1005 9.2.2. Metric for Efficiency

1006 We use *throughput* as the evaluation metric for efficiency. Specifically, throughput is defined as the total number of rollout  
 1007 environment frames divided by the total wall-clock time of one training epoch, which approximately equals the sum of rollout  
 1008 time and training time.

### 1009 9.2.3. Tasks for Efficiency

1010 We evaluate our framework on all tasks in ManiSkill, LIBERO and RoboTwin.

1011 For **ManiSkill**, we select the “PutCarrotOnPlateInScene-v2” task for a quick test. It is adapted from  
 1012 “PutCarrotOnPlateInScene-v1” in ManiSkill, with modifications to the reset function. In the original reset implementa-  
 1013 tion, additional simulation steps were applied to ensure that all objects remained static. However, this significantly increased  
 1014 the reset time compared to the step function and was incompatible with partial reset. In practice, we observed that these extra  
 1015 steps had a negligible impact on object states in the “PutCarrotOnPlateInScene” task, since the initial position of the carrot  
 1016 was already carefully calibrated. We therefore delete the additional simulation steps in the reset function. The evaluation is  
 1017 conducted on 8, 16, and 32 NVIDIA H100 (80GB) GPUs, and we use 256 parallel environments, each running for 80 steps.

1018 For **LIBERO**, we adopt the LIBERO-Long task set for a quick test, which is the longest task set in LIBERO. The number  
 1019 of parallel environments is set to 64, 128, and 256 for 8-, 16-, and 32-GPU setups, respectively, with the corresponding  
 1020 number of environment steps set to 4096, 2048, and 1024.<sup>3</sup>

1021 For **RoboTwin**, we select the `place_empty_cup` task for a quick test, which requires the robot to use an arm to place  
 1022 the empty cup on the coaster. The number of parallel environments is set to 64, 128, and 256 for 8-, 16-, and 32-GPU setups,  
 1023 respectively, with the corresponding number of environment steps set to 800, 400, and 200. We have the same GPU allocation  
 1024 setting as LIBERO-Long, due to the GPU resources being the upper bound for the rollout.

### 1025 9.2.4. Baselines for Efficiency Experiments

1026 For ManiSkill, no existing framework supports multi-GPU rollout and training. Thus, we take the naive disaggregated  
 1027 allocation mode as the baseline, and compare it with our colocated mode and hybrid modes, where the hybrid mode includes  
 1028 both one-stage and two-stage fine-grained pipelining configurations.

1029 For LIBERO and RoboTwin, we compare against SimpleVLA-RL [20], an open-source framework for RL of VLA Models  
 1030 training built on VeRL. Since SimpleVLA-RL only supports the colocated mode, we use it as the baseline and additionally  
 1031 evaluate our colocated mode and hybrid mode (with one-stage and two-stage fine-grained pipelining). We omit the disaggre-  
 1032 gated mode results for LIBERO and RoboTwin, as during training the *Training* component never reuses the same GPUs as the  
 1033 rollout process. Consequently, even when accounting for the offload-onload overhead, the hybrid (one-stage) configuration  
 1034 provides superior performance to the disaggregated mode.

## 1035 10. Additional Experimental Results

### 1036 10.1. Evaluation of OOD for ManiSkill

1037 Table 6 is the detailed results ManiSkill’s experiments.

<sup>3</sup>Since vectorized environments in LIBERO rely on multi-processing, the number of physical CPU cores per node becomes the upper bound for efficient rollout. Therefore, we scale the number of parallel environments in proportion to the number of nodes, and consequently to the total number of GPUs.

Table 6. Detailed breakdown of Out-Of-Distribution evaluation results on ManiSkill. Values denote success rates (%).

Method	Out-Of-Distribution		
	Vision	Semantic	Execution
OpenVLA (Base)	38.75	35.94	42.11
OpenVLA (RLinf-GRPO)	74.69	72.99	77.86
OpenVLA (RLinf-PPO)	82.03	<b>78.35</b>	<b>85.42</b>
OpenVLA-OFT (Base)	27.73	12.95	11.72
OpenVLA-OFT (RLinf-GRPO)	84.69	45.54	44.66
OpenVLA-OFT (RLinf-PPO)	<b>92.11</b>	64.84	73.57

## 10.2. Ablation Studies

1038

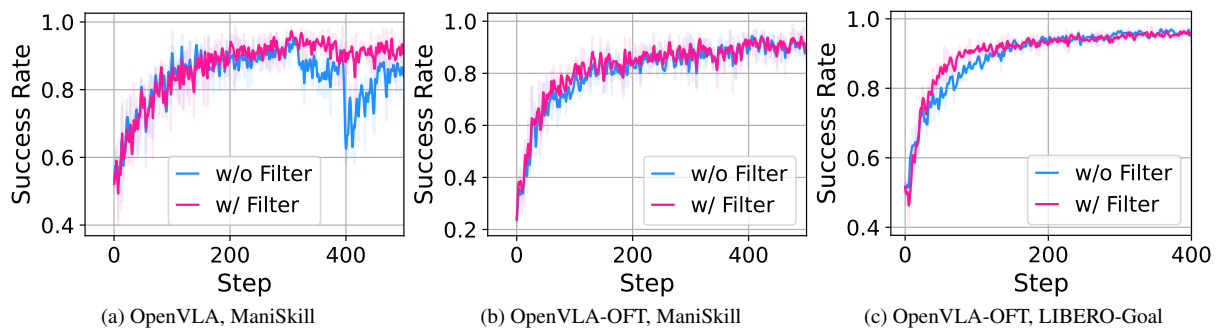


Figure 15. Ablation studies on success rate filtering under different settings.

**Success rate filtering** As shown in Figure 15, we implement a success rate filter for GRPO that discards groups in which all trajectories have identical cumulative rewards. This mechanism can improve the stability of GRPO training. For instance, in the OpenVLA ManiSkill setting, training without the filter (“w/o Filter”) exhibits a clear collapse around step 400, whereas enabling the filter (“w/ Filter”) largely alleviates this issue. However, the benefit of the filter is not universal: in the OpenVLA-OFT ManiSkill setting and the OpenVLA-OFT LIBERO-Goal setting, its effectiveness is much less pronounced.

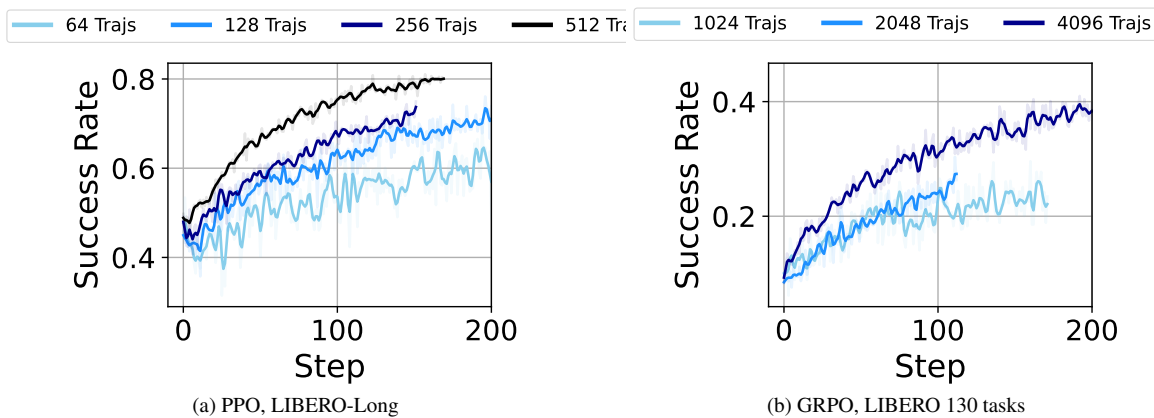
1039  
1040  
1041  
1042  
1043

Figure 16. Ablation study on rollout data size. Darker colors correspond to larger rollout datasets.

1044 **Effect of rollout data size** The rollout batch size has a notable influence on RL performance, especially for on-policy  
 1045 algorithms such as PPO and GRPO. In general, larger rollouts per epoch enable more substantial policy improvement within  
 1046 each training iteration. As illustrated in Figure 16, larger rollouts consistently achieve higher success rates when evaluated  
 1047 by training epochs.

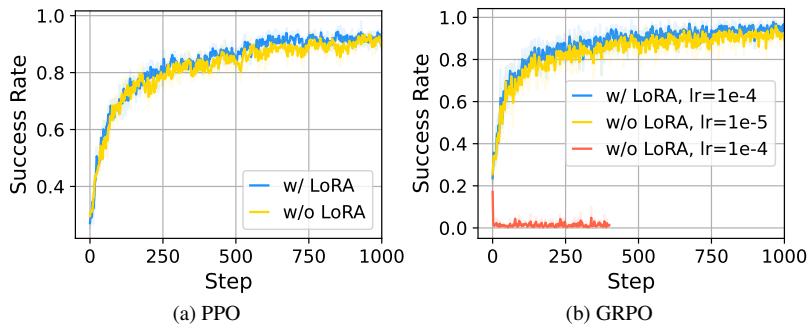


Figure 17. Ablation of LoRA in ManiSkill with OpenVLA-OFT.

1048 **Using LoRA may not directly affect performance but often requires different hyperparameters** Figure 17 shows the  
 1049 training curves with and without LoRA [13]. The x-axis denotes the number of training epochs, while the y-axis shows  
 1050 the success rate in ManiSkill. The curves are overall similar, suggesting that LoRA itself does not substantially change  
 1051 performance. However, the choice of whether to use LoRA can influence the optimal hyperparameters. For example, in  
 1052 GRPO experiments (Figure 17b), using the same learning rate of  $1 \times 10^{-4}$  leads to normal improvement with LoRA, but  
 1053 the success rate without LoRA collapses to zero. In contrast, when the learning rate is reduced to  $1 \times 10^{-5}$ , the non-LoRA  
 1054 setting also achieves stable improvement. These results indicate that different LoRA configurations may require separate  
 1055 hyperparameter tuning.