Refined Policy Distillation: From VLA Generalists to RL Experts

Tobias Jülg¹, Wolfram Burgard¹ and Florian Walter¹

Abstract—Recent generalist Vision-Language-Action Models (VLAs) can perform a variety of tasks on real robots with remarkable generalization capabilities. However, reported success rates are often not on par with those of expert policies. Moreover, VLAs usually do not work out of the box and often must be fine-tuned as they are sensitive to setup changes. In this work, we present Refined Policy Distillation (RPD), an RLbased policy refinement method that enables the distillation of large generalist models into small, high-performing expert policies. The student policy is guided during the RL exploration by actions of a teacher VLA for increased sample efficiency and faster convergence. Different from previous work that focuses on applying VLAs to real-world experiments, we create finetuned versions of Octo and OpenVLA for ManiSkill2 to evaluate RPD in simulation. As our results for different manipulation tasks demonstrate, RPD enables the RL agent to learn expert policies that surpass the teacher's performance in both dense and sparse reward settings. Our approach is even robust to changes in the camera perspective and can generalize to task variations that the underlying VLA cannot solve.

I. Introduction

Recently, robot learning has seen a paradigm shift from Reinforcement Learning (RL) of task-specific policies to large generalist Vision-Language-Action Models (VLAs) that are trained in a supervised fashion using imitation learning techniques [1], [2], [3], [4], [5], [6], [7], [8]. VLAs are easy to train and show good performance on different realworld manipulation tasks with a single model. However, they usually need to be trained on a large number of human expert demonstrations, limiting their performance to the amount and variety of the training data. Different from the field of natural language processing, creating more data is very costintensive in robotics as it usually means collecting trajectories by humans through time-consuming teleoperation. Due to limited training data, VLAs have problems generalizing across setups [9]. Applying VLAs in new environments, thus, often requires fine-tuning on domain-specific data, requiring human effort yet again to record demonstrations, which further limits their widespread adaptation. Only recently, Xu et al. [10] utilized data generated by RL agents for VLA fine-tuning, which can significantly scale up setup-specific fine-tuning.

RL can, in principle, scale much better than imitation learning as an agent learns to maximize rewards through autonomous interaction with its environment and, thus, generates its own data. Often, training is performed in simulation, which can be massively parallelized [11], [12], [13].

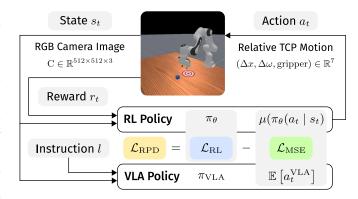


Fig. 1. The architecture of RPD: We distill a VLA policy into an RL student policy. This has two effects: First, the RL agent is bootstrapped with guided exploration through the VLA teacher. Second, the student policy can interact with the environment and is, thus, able to surpass the VLA's performance, refining the distilled policy.

RL policies are also considerably more compact than VLAs as they have fewer parameters, resulting in faster inference and increased efficiency.

However, applying RL in robotics also poses several challenges: Training is very sensitive to the hyperparameter configuration and often needs environment-specific tuning [14]. RL agents require extensive exploration of the state space, which is usually sample-inefficient and precludes execution on a real robot. Simulated environments are, therefore, essential for most RL algorithms to work. The resulting sim-to-real gap can be closed by subsequent fine-tuning on the real robot.

In this work, we introduce *Refined Policy Distillation* (RPD) as a novel method for distilling and optimizing expert policies from generalist VLA policies. The idea behind RPD is that a VLA policy trained on large-scale datasets stores general task knowledge that can inform the training process of an RL algorithm, even if the policy by itself can solve the task only at a low or possibly even zero success rate. Our evaluation shows that RPD is not only capable of distilling and improving policies from finetuned VLA policies but that it can also successfully distill VLA knowledge while learning unseen tasks. This also applies when the camera angle changes, a situation in which current VLAs largely fail. RPD extends the Proximal Policy Optimization (PPO) [15] algorithm and only requires a slight modification of its loss function.

The main contributions of this work can be summarized as follows:

 We propose RPD, a novel method for distilling and refining large VLAs with on-policy RL into compact expert policies.

¹All authors are with the Department of Computer Science and Artificial Intelligence, University of Technology Nuremberg, Germany. Contact: [tobias.juelg, wolfram.burgard, florian.walter]@utn.de

- We demonstrate RPD's increased sample efficiency and improved stability on six different in-distribution tasks for both dense and sparse reward settings.
- We show its positive generalization on two out-ofdistribution task variations and a changed camera perspective.

II. RELATED WORK

In recent years, several approaches have been developed for learning generalist robot policies conditioned on vision and language inputs. One of the most common ways is to use an existing VLM backbone trained on large amounts of internet-scale data and fine-tune it on recorded robot trajectories [3], [5], [7], [8], most notably from the Open X-Embodiment dataset [16]. Other models are based on large transformers and trained on vision and language or goalimage instructions from scratch [1], [2], [4], [17], [18]. We refer to both of these transformer-based policy types as VLAs.

Not all VLAs are publicly available [3], and for some, only the weights have been released, but not the source code [2]. This precludes fine-tuning them on new tasks or setups. However, both the code and the weights of many recent models are openly available [4], [5], [7], [17]. In this work, we use Octo [4] and OpenVLA [7] for our tests as they represent two important classes of models, but RPD works with any other VLA.

Policy Distillation (PD) is a widely investigated method in RL to distill teacher policies into student policies to increase training speed and decrease the complexity of the resulting policies [19]. Early work in this domain investigated the distillation of DQN agents [20]. The authors compared different approaches, including KL divergence and MSE loss between the Q functions. KL divergence was used to distill DQN policies into PPO policies, too [21]. PD can also be framed as a supervised policy transfer problem that allows the distillation from a teacher that uses a different RL algorithm than the student [22]. Compared to our method, those works only distill from RL teacher policies but not from generalist policies.

A recent work in this domain is Proximal Policy Distillation (PPD) [23], which distills a PPO teacher policy [15] into a potentially larger PPO student policy. PPD only samples from the student distribution, which allows it to lower the exploration bias induced by the teacher. To incorporate the teacher's policy, PPD modifies the PPO loss by adding the KL divergence between the teacher's and the student's policy distribution together with a clipping term. This method is transferable to VLA teacher policies as long as they are stochastic. However, sampling the action distributions from a VLA, which is required for estimating the KL divergence, is compute-intensive and multiplies GPU load and memory usage by the sample size. In this work, we employ a simple MSE loss, which avoids compute-intensive sampling of actions from a large-scale VLA.

Behavior Cloning (BC) is an umbrella term for offline imitation learning: Given a dataset of state transitions, the

goal is to learn a policy that follows the trajectories from the dataset as closely as possible [24]. Exact formulations differ, but the main idea is to maximize the likelihood of the state transitions in the dataset. BC was also combined with RL to speed up the agent's exploration through a fixed prerecorded dataset [25], [26]. Recently, the term Offline RL has been used to describe similar approaches but they focus mainly on off-policy RL methods [27], [28]. Offline RL is related to PD as both approaches aim to learn from expert knowledge but uses datasets instead of learned policies. We employ methods similar to BC to process external actions in the RL training loop but leverage pre-trained generalist policies to generate them.

Xu et al. [10] investigated dataset generation with an offline soft actor-critic variant [28], [29] for VLA fine-tuning. While their experiments show performance improvements of the VLA policy of up to 40%, the method still requires already trained RL expert policies. By contrast, our method distills large foundation models into smaller on-policy RL agents, which avoids training RL policies from scratch. Through interaction with the environment, these policies can surpass the performance of their VLA teachers on specific tasks and become expert policies.

In summary, compared to these previous works, our approach distills from large VLA models using RL. The resulting policies are compact and fast expert policies. In all of our experiments, the resulting performance exceeds the VLA's performance, which we refer to as *policy refinement*.

III. METHODOLOGY

RPD is a method for distilling large generalist VLA policies into small and fast task-specific expert policies. Standard supervised BC techniques usually learn policies that do not exceed the performance of the teacher For this reason, we combine policy distillation with RL. This helps the student agent to explore more efficiently and reach higher rewards faster, even in sparse reward settings. In addition, it also allows for refinement of the teacher policy to surpass its performance and become an expert policy for the distilled task. At the core of RPD is a PPO RL agent with a modified objective function that includes a Mean Squared Error (MSE) term between the action mean predicted by PPO and the expectation of the VLA action to pull the action mean of the RL policy closer to the predictions of the VLA. Fig. 1 provides an overview of how RPD integrates VLA actions into the agent-environment RL training loop.

A. Refined Policy Distillation

RPD combines RL with a Behavior Cloning (BC) objective. For RL, we model the interaction between an agent and its environment as a finite Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, p, \rho_0, \gamma)$ [30]. \mathcal{S} is the state space, \mathcal{A} the action space, $r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ the reward function, $p(s_{t+1} \mid s_t, a_t)$ the state transition probability distribution of the environment, ρ_0 the probability distribution over the initial states s_0 , and $\gamma \in [0,1]$ the discount factor. Every task has its own MDP associated with it. We consider stochastic

policies $\pi_{\theta}(a_t \mid s_t)$ that are parameterized by weights θ . The goal of RL is to find a set of weights θ^* that maximizes the expected sum of discounted rewards:

$$\theta^* = \arg\max_{\theta} \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right]$$
 (1)

with $a_t \sim \pi_{\theta}(a_t \mid s_t)$, $s_0 \sim \rho_0$, $s_{t+1} \sim p(s_{t+1} \mid s_t, a_t)$, and the episode length T.

BC trains a policy from expert demonstrations that are usually provided as a dataset $\rho_{\mathcal{D}}$ of state transitions. We adopt a BC objective [25] that maximizes the likelihood of the policy:

$$\max_{\theta} \sum_{(s,a) \in \rho_{\mathcal{D}}} \ln \pi_{\theta}(a \mid s). \tag{2}$$

The transitions can also be sampled online and do not need to be stored in a dataset. In our case they come from the VLA policy π_{VLA} . Under the assumption of Gaussian action distributions, the BC objective essentially boils down to an MSE minimization objective between the two distributions' means and a term for the variance that we assume to be constant and therefore neglect.

The RPD objective combines PPO with BC through a modified objective function:

$$\mathcal{L}_{RPD}(\theta) = \mathcal{L}_{RL}(\theta) - \mathcal{L}_{MSE}(\theta)$$
 (3)

where \mathcal{L}_{RL} is the standard PPO clipped surrogate objective function

$$\mathcal{L}_{RL}(\theta) = \mathbb{E}_t \left[\min(r_t(\theta) \hat{A}_t, c(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$
(4)

with clip function c and the notation from Schulman et al. [15]. \mathcal{L}_{MSE} is the expectation of the MSE between the action mean $\mu(\pi_{\theta}(a_t \mid s_t))$ predicted by the PPO policy and the expectation $\mathbb{E}\left[a_t^{\text{VLA}}\right]$ of the actions $a_t^{\text{VLA}} \sim \pi_{\text{VLA}}(a_t \mid s_t)$ returned by the VLA policy:

$$\mathcal{L}_{\text{MSE}}(\theta) = \mathbb{E}_t \left[\left(\mu(\pi_{\theta}(a_t \mid s_t)) - \mathbb{E}\left[a_t^{\text{VLA}} \right] \right)^2 \right]$$
 (5)

Note that the PPO objective is maximized while the MSE is minimized, which is why it appears with a negative sign. In practice, we found that sampling a single VLA action per step already yields good performance. Increasing the sample size would slow down the training considerably and, therefore, was not further investigated. Our training objective is to find a set of optimal weights θ^* that maximize the objective function \mathcal{L}_{RPD} :

$$\theta^* = \arg\max_{\theta} \mathcal{L}_{RPD}(\theta)$$
 (6)

In addition to the MSE-based RPD objective, which we refer to as RPD-MSE, we will also consider variants where the MSE loss is replaced by other loss functions.

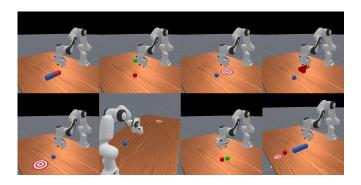


Fig. 2. Overview of the eight different tasks from ManiSkill2 that were used to distill expert policies with RPD. See Tab. I for the task names (tasks are depicted in row-major order).

B. Tasks, Dataset, and VLA Fine-Tuning

We use eight different manipulation tasks from ManiSkill2 [31] to evaluate RPD. Fig. 2 provides a visual overview. Our criteria for choosing these tasks were the availability of demonstrations that are no longer than the maximum episode length for reasons detailed below and that all action spaces contain relative task space actions a with the following format:

$$a = (\Delta x, \Delta \omega, \text{gripper}) \in \mathbb{R}^7$$
 (7)

The camera perspectives shown in Fig. 2 match the environment state observed by both the VLA and the RL agent. Moreover, the same camera was used to record the VLA fine-tuning data. ManiSkill2 calls this the human camera. It is normally not meant to be used for training as it contains visual hints such as a green sphere for the goal location of the cube in the PickCube task.

Goal locations in ManiSkill2 are drawn randomly and provided to the agent only through explicit state information. The VLAs considered in this work, however, only receive RGB camera state information and a fairly simple natural language instruction. Therefore, at least one of these two input modalities needs to provide a hint for the currently selected target location. While encoding the target in the language instruction appears to be a natural choice, the precision may be limited as VLAs are usually pre-trained on simple language instructions that do not contain exact spatial information. We leverage the visual cues in ManiSkill2's human camera to encode all relevant task information. Another reason for using this camera is that its perspective shows the scene from the side rather than from the front. This side view corresponds better to the camera perspective used in many setups of the Open X-Embodiment dataset.

Initially, we evaluated the published checkpoints for OpenVLA and Octo on the selected ManiSkill2 tasks using the human camera and the language instructions shown in Tab. I. However, as reported in the table, both models had a success rate of zero on all tested tasks. This real-to-sim gap is expected, as both Octo and OpenVLA were only trained on real-world datasets. When using ManiSkill2's raytracing renderer, the success rate of Octo increased to

TABLE I: EVALUATION REWARDS (RW) AND SUCCESS RATES (SR) OF THE VLAS ON THE SELECTED MANISKILL2 TASKS.

Task ID	Language Instruction	Fine-tuned	openvl RW	a-base SR	openvl RW	a-finetuned SR	octo-	base SR	octo-fi RW	netuned SR
	II	<u> </u>								
LiftPegUpright-v1	"lift the peg upright"	✓	0.20	0%	0.46	12%	0.20	0%	0.41	0%
PickCube-v1	"pick up the cube"	✓	0.06	0%	0.18	6%	0.06	0%	0.17	9%
PullCube-v1	"pull the cube towards the robot base"	✓	0.07	0%	0.27	92%	0.07	0%	0.24	90%
PullCubeTool-v1	"pull the cube by using the red tool"	×	0.10	0%	0.04	12%	0.10	0%	0.03	10%
PushCube-v1	"push the cube away from the robot base"	✓	0.10	0%	0.12	27%	0.10	0%	0.16	67%
RollBall-v1	"push the ball"	✓	0.02	0%	0.03	4%	0.02	0%	0.06	10%
StackCube-v1	"stack the red cube on the green cube"	✓	0.05	0%	0.13	3%	0.05	0%	0.13	1%
PokeCube-v1	"push the cube by using the blue tool"	×	0.06	0%	0.09	12%	0.06	0%	0.08	12%

more than 2% on PushCube, indicating that some positive domain transfer is possible for more realistic-looking input data. Still, those success rates are too small for a teacher policy, and ManiSkill2 does not support GPU rendering of the human camera with raytracing, which is required for RL.

In order to evaluate RPD under realistic conditions, we first fine-tuned Octo and OpenVLA on the simulated environment using RL-generated expert demonstrations provided by ManiSkill2. All demonstrations were rendered with the human camera at an image resolution of 256×256 . During training, we found that increasing the maximum episode length to 300 steps led to considerably higher success rates as the ManiSkill2 dataset contains a subset with episodes longer than the maximum of the 50 steps used for most environments. As long episodes can make RL algorithms behave less stable and also require a higher number of VLA inference steps, we excluded them from fine-tuning, which led to the expected success rates for both Octo and OpenVLA. The results and language instructions for all tasks and models are summarized in Tab. I.

The two tasks PullCubeTool and PokeCube are excluded from the fine-tuning dataset to test cross-task generalization. While the normalized rewards are low, as expected, the reported success rates seem unreasonably high. This is because in these two tasks, the agent is supposed to grasp a tool and use it to either pull or poke a cube. We found that the VLAs were just pulling or pushing the cube to the target location without utilizing the tool, which is counted as a success by the sparse reward implemented by ManiSkill2 as it just compares the distance towards the goal location and not whether the tool was actually used. The dense reward, however, takes tool use into account and accurately shows the decreased performance. Nevertheless, we still decided to keep both tasks. Even if the learned behavior is just pushing and pulling, which the success rate can still accurately display, it still differs from PushCube and PullCube because of distracting objects in the scene. The fact that the language instruction is partly ignored by the VLA in that case is not relevant for testing the generalization behavior of RPD.

C. VLA Integration Details

To integrate different VLAs into the RL training loop, we defined a generic policy interface that takes a batch of visual observations and language instructions as input and returns a relative task space motion command including the

gripper state. A technically challenging part is the software dependency management for the implementations of the different VLAs, as they are partly incompatible with each other and with the dependencies of the RL training code and the simulator. For example, Octo requires a specific version of JAX, which depends on a CUDA version that does not work with the CUDA version that OpenVLA's PyTorch version needs.

To solve this issue, we created separate Python environments, each of which contains the correct dependencies. The generic VLA interface is exposed via an HTTP API server that runs with the correct environment. During training, the RL agent connects to this server and sends batches of observations.

A new problem that comes with this approach is serialization. Communication can quickly become a bottleneck when large batches of raw image data have to be serialized. If we run both VLA inference and the RL training on the same machine, we use shared memory to avoid serialization.

Lastly, OpenVLA does not support batch processing in its inference code. This limitation further slows down the already large model, as it complicates the implementation of parallel inference. Moreover, the CPU-based preprocessing occupies excessive CPU resources, necessitating manual code optimizations to move it to the GPU. Even after utilizing a parallelization factor of eight (running two OpenVLA instances on each GPU within a four Nvidia H100 node setup), it remains considerably slower than Octo operating on a single Nvidia A40 GPU. Due to these constraints, we were only able to conduct a subset of our experiments using OpenVLA.

IV. EXPERIMENTAL RESULTS

We evaluated RPD on two VLAs, namely Octo [4] and OpenVLA [7]. The latter model with its 7B parameters is much larger than the former with only 93M parameters and therefore requires also considerably more computing time, which is exacerbated by the lack of support for batched inference. This is why we provide evaluation results for OpenVLA only for a subset of the tasks. henever the base model is not explicitly mentioned we refer to Octo.

We split the evaluation into four parts: First, we test RPD on the challenging PickCube task to explore the performance differences between loss variants and baselines. We then evaluate RPD and its PPO baseline on a total of six different

tasks that are included in the VLA fine-tuning dataset. RL often struggles with sparse rewards due to inefficient exploration. We argue that RPD guides this exploration process and, thus, should also improve sample efficiency and performance on sparse rewards. To test this hypothesis we also run the tasks mentioned above in a sparse reward setting. Finally, we evaluate RPD's cross-task and cross-setup generalization capabilities by testing it on two tasks that it was not trained on and a variant of the PushCube task with an altered camera perspective. The latter simulates small setup changes that can easily happen in reality and have a considerable effect on the performance of current VLMs.

A. RPD Variants and Baselines

We compare three different variants of RPD on ManiSkill2's Cube task:

- 1) RPD-MSE, the main RPD variant as defined in (3).
- 2) RPD-L1, which uses an L1 loss instead of the MSE.
- 3) RPD-BC, which uses a naive maximum likelihood loss as defined in (2).

All RPD variants are trained with PPO and constant hyperparameters which were adopted without modification from the ManiSkill2 baseline [31]. The only exception is the batch size that had to be adjusted to fit the VLA on the same GPU as the RL training loop. Due to compute constraints we distilled from OpenVLA only with RPD-MSE.

We train policies with vanilla PPO and PPD as baselines and use the PPO implementation from CleanRL [32] that is included in ManiSkill2. As there is no source code publicly available for PPD, we use our own implementation, which extends CleanRL's PPO. To compute the KL divergence PPD requires for its distillation loss, we sample ten actions from the VLA and fit a multivariate Gaussian distribution to it. All hyperparameters are set to ManiSkill2's defaults. The weighting factor of PPD is set to $\lambda=1.0$.

Fig. 3 depicts the validation success rates during training for all RPD loss variants and baselines mentioned above. For reference, it also includes the performance of the VLAs after fine-tuning on the ManiSkill2 dataset previously presented in Tab. I.

The evaluation shows that RPD-MSE outperforms all other variants, closely followed by RPD-L1. The distilled policy quickly surpasses Octo and OpenVLA to become a refined policy. RPD-MSE with OpenVLA performs slightly worse than the distillation with Octo. The PPO baseline also converges to similar success rates as RPD-MSE and RPD-L1 of around 80% but learns substantially slower and shows much larger fluctuations in the training process. RPD-BC slightly surpasses Octo's performance threshold but stays at that level and does not refine the policy. Lastly, PPD fails to lift the policy above the Octo baseline.

The reason why our PPD baseline does not manage to learn the task successfully is likely that it requires the teacher and student policy distributions to be of similar shape for the KL divergence to converge. However, we found that the action distributions of the VLA policies are sometimes bimodal, which results in non-optimal parameters for the

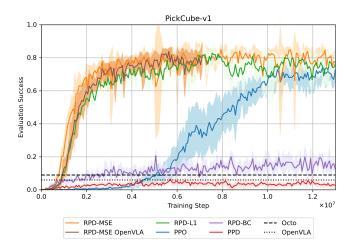


Fig. 3. Average success rates during training for vanilla PPO, PPD, and three different RPD variants: RPD-MSE, RPD-L1, and RPD-BC, which all distill from Octo over five runs with different seeds. indicate standard deviations. We also evaluated RPD-MSE on OpenVLA. Due to computing constraints, we could only perform a single training run. The performance of the fine-tuned VLAs is indicated by dashed lines.

approximated normal distributions that we use to compute the KL divergence in PPD. We also speculate that in our setup PPD depends strongly on the teacher's performance. Lastly, PPD might also be sensitive to PPO hyperparameters, which we have not investigated. We did, however, change the parameter λ , which did not lead to substantially different results.

Even though RPD-BC uses a very similar loss compared to RPD-MSE, it includes a variance component that forces the standard deviation to smaller and smaller values. This leads to overfitting of the teacher's policy and, thus, does not converge to a considerably better performance. RPD-BC can therefore be seen as a policy distillation method without refinement. It may be useful for situations where a close match with the teacher policy is required.

RPD-MSE likely performs better on Octo than on OpenVLA because of the performance difference already present in the foundation models themselves. Both RPD-MSE and RPD-L1 show quicker and less noisy convergence than the PPO baseline with the same hyperparameters. This means that RPD efficiently guides the rather inefficient exploration process of the RL agent, leading to larger rewards in less time, which improves sample efficiency and robustness.

B. RPD on Different Tasks

Fig. 4 shows the training curves for RPD-MSE with Octo on six different manipulation tasks that were part of the training dataset. In all six cases, RPD learns faster than the PPO baseline and often converges to higher success rates. This increased performance can also be seen in the reward plots. For RollBall and StackCube, the ManiSkill2 PPO hyperparameters fail to solve the tasks whereas RPD converges at least to partial success using the very same parameters. The decreased standard deviation also shows that

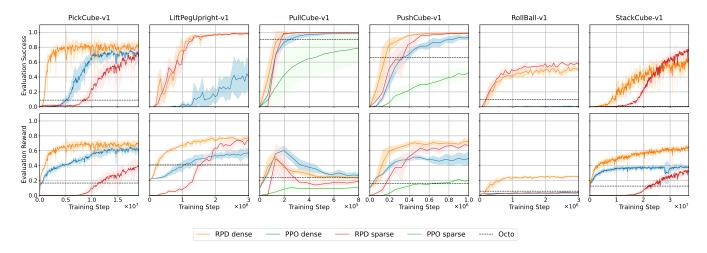


Fig. 4. Success and reward curves of RPD with Octo and vanilla PPO for both dense and sparse rewards on the six different ManiSkill2 tasks tha part of the fine-tuning dataset. The values are averaged over five runs with different seeds and are recorded in an evaluation environment. Standard deviations are indicated by shaded areas. All training runs were performed with the same hyperparameters taken from the ManiSkill2 PPO baseline. For all tasks, RPD outperforms Octo quickly and converges faster than vanilla PPO. In some cases, it even finds good policies when vanilla PPO fails. This effect increases for sparse rewards.

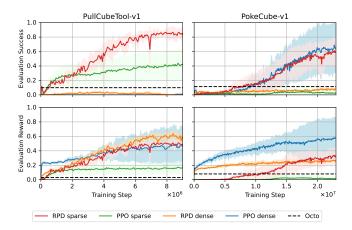


Fig. 5. Success and reward curves for RPD applied to Octo and vanilla PPO on tasks that are *not* part of Octo's fine-tuning dataset. The curves show average results for five training runs on different seeds and are recorded in evaluation environments. The shaded areas indicate standard deviations. Note that the Octo baseline does not represent the correct reward in the dense reward setting as it solves the tasks without tool use. See sec. III for details.

RPD is generally more robust than the vanilla PPO baseline. These results demonstrate that RPD improves PPO's training speed consistently and over many tasks while surpassing its final performance. It also refines the distilled policy on dense rewards.

C. Sparse Rewards and Sample Efficiency

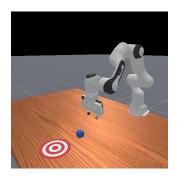
ManiSkill2 defines its sparse rewards as follows: -1 for failed cases, 1 for success, and 0 otherwise. This means that non-zero reward is only given at the end of the episode, which is different from dense rewards where every action usually receives a non-zero reward that is shaped to guide the agent towards a certain behavior. Sparse reward problems are often a big challenge for RL algorithms, as relevant

information about the agent's performance is only obtained once the task has been solved successfully at least one time. Thus, in the beginning, the agent needs to complete an episode successfully just through random exploration, which makes the training process very inefficient compared to RL with dense rewards. As a result, it can occur that an agent learns a task perfectly fine with dense rewards but fails in the corresponding sparse reward setting, as can be seen for PPO in Fig. 4 for some of the tasks.

Fig. 4 also shows training runs with the six selected ManiSkill2 tasks for sparse rewards. Note that the resulting reward values cannot be directly compared to the dense reward setting. The only modification we applied to the hyperparameters from the ManiSkill2 baselines was to set γ from 0.8 to 0.99, which avoids vanishing returns for early episode steps. In all six cases, RPD on sparse rewards consistently outperforms vanilla PPO. In most tasks, the baseline completely fails to learn the task, whereas RPD succeeds. It makes sense that RPD can realize its full potential on sparse rewards as it is even more important in this setting to have a teacher policy that enables the student to successfully complete episodes and, thus, retrieve environment rewards that drive the learning process. In summary, RPD drastically improves sample efficiency for on-policy RL in both sparse and dense reward cases.

D. Generalization to New Environments

We also evaluated RPD with Octo on tasks that were held out of the VLA fine-tuning dataset to test its generalization capabilities. Ultimately, RPD is limited by the performance of the teacher and, in this case, by the cross-task generalization of the VLA. Still, it is worth investigating to what extent the degraded VLA performance will influence the RPD training. We study this question by looking at task variations, which means they are out-of-distribution with respect to the fine-tuning dataset. More specifically, we evaluate



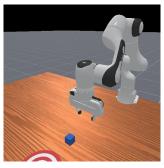


Fig. 6. The PushCube task with a variation of the camera perspective. The left side shows the perspective used for fine-tuning Octo. The right side depicts the perspective we used to test RPD's performance in situations where the VLA needs to generalize from minor setup changes, a task where state-of-the-art VLAs fail.

RPD's performance on the hold-out tasks PullCubeTool and PokeCube as well as on a variation of the PushCube task with a different camera perspective.

a) Hold-Out Tasks: As mentioned in sec. III, generalization to unseen tasks can be measured based on the performance on the hold-out tasks (see Tab. I). However, it requires a reinterpretation: PullCubeTool and PokeCube have to be interpreted as PullCube and PushCube with distracting objects. For the original versions that require tool usage, we could not observe high generalization capabilities.

The results of our experiments are summarized in Fig. 5. On the one hand, RPD fails to learn the hold-out tasks on dense rewards as they force the RL agent to use the tool and the VLA policy consequently cannot provide any meaningful actions. On the other hand, in the sparse reward setting, the RL agent learns the reinterpreted tasks, for which the VLA policy generalizes well. Therefore, it makes sense that RPD learns these tasks successfully as shown in Fig. 5. It outperforms PPO on both tasks even though the reduced task difficulty is the same for both algorithms.

b) Change in Camera Perspective: VLAs have traditionally been very sensitive to camera adjustments, leading to severe performance degradation [9]. To counteract performance drops VLAs often require fine-tuning with camera data from the new angle, camera model, etc. This problem also occurs with our fine-tuned versions of Octo and OpenVLA: When we only slightly change the camera angle in the PushCube task, as shown in Fig. 6, the success rates of OpenVLA drops from 27% to 4.5%. Octo even drops from 67% to 0%. The latter is likely a result of the less powerful vision encoder that is utilized in Octo which may overfit the training data.

Fig. 7 shows the success rates of vanilla PPO and RPD-MSE for both Octo and OpenVLA on the changed camera perspective from Fig. 6 on the PushCube task. PPO performs similar on our task variation as in the original PushCube task with sparse rewards, see Fig. 4. This is expected because the agent's visual encoder is learned from scratch in the RL training. Accordingly, there is no bias towards a specific camera angle as long as the provided

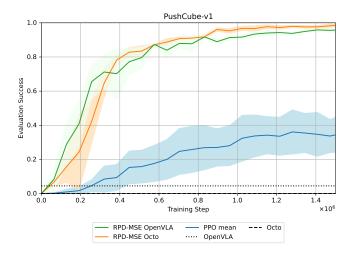


Fig. 7. Average evaluation success rates of vanilla PPO and RPD-MSE with both Octo and OpenVLA on a variation of the PushCube task with sparse rewards over five different seeds. Shaded areas indicate standard deviations. The camera perspective was changed as shown in Fig. 6 to investigate the distillation performance of the degraded VLAs.

information is sufficient to solve the task. Both RPD variants manage to outperform PPO.

Similar to the previous paragraph, where distracting objects are added to the scene, RPD is robust against visual changes despite a considerable drop of the VLA's performance. The reason why RPD can benefit from degraded VLA actions is because they still guide the exploration in a good direction. Importantly, low success rates do not necessarily mean that the agent's actions are bad. For example, the VLA agent may almost complete the task but then perform a wrong action in the end. In this case, the success rate is zero although the actions are still valuable for a student policy, as can be seen in the evaluation results for RPD.

E. Limitations

The main limitation of RDP is the performance of the VLA on the task at hand. If the VLA cannot provide any meaningful actions due to domain shift or because the task was not included in the training dataset, RPD can be seen as a new exploration component in the loop. Depending on the resulting action distribution, the VLA may provide bias that is useful for learning the task or could improve the exploration of the environment.

However, if the actions generated by the VLA cannot provide such helpful bias, RPD may also decrease the RL agent's performance. As reported for the training with dense rewards on the hold-out tasks, if the task and no variation of it have ever been seen by the VLA, its actions hinder the PPO training in RPD, and vanilla PPO performs better. hen PPO is already tuned to a task or the task is very easy, PPO is already sample-efficient and RPD may only provide limited improvement.

V. CONCLUSION AND FUTURE WORK

This work introduced RPD, a novel on-policy RL-based approach for distilling generalist VLAs into refined expert

policies. To the best of our knowledge, our work is the first that distills VLAs into compact expert policies. Our experiments show that RPD not only makes the RL training more efficient by guiding its exploration through the actions of the VLA but also that it can leverage task knowledge from the VLA even in situations where the generalist policy fails, such as when the camera perspective changes. Thus, RPD can help to improve the usability and accessibility of VLAs.

Despite the reported efficiency gains in the RL process, RPD still requires training in simulation, which is why we fine-tuned two state-of-the-art generalist VLA policies, Octo and OpenVLA, on simulation datasets. As a side contribution, our work shows that VLAs can be successfully utilized in experiments in simulated environments. This provides a promising perspective for scientists who have no access to a physical robot setup to conduct research on VLAs. An interesting direction for future research is to combine RPD with sim-to-real methods [33] for deployment on real robots.

Another promising line of research is to use the distilled expert policies for collecting new training data. Following the method proposed by Xu *et al.* [10] to fine-tune VLAs, it will be possible to improve generalist policies without collecting any additional human demonstration data.

ACKNOWLEDGEMENT

This work is supported by the project GeniusRobot funded by the German Federal Ministry of Education and Research (BMBF grant no. 01IS24083). The authors gratefully acknowledge the scientific support and HPC resources provided by the Erlangen National High Performance Computing Center (NHR@FAU) of the Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU) under the BayernKI project no. v108be. BayernKI funding is provided by Bavarian state authorities.

REFERENCES

- [1] S. Reed, K. Zolna, E. Parisotto, et al., "A generalist agent," Proc. of the Int. Conf. on Learning Representations (ICLR), 2024.
- [2] A. Brohan, N. Brown, J. Carbajal, et al., "RT-1: Robotics transformer for real-world control at scale," in *Proc. of Robotics: Science and Systems (RSS)*, 2023.
- [3] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, and K. Choromanski, "RT-2: Vision-language-action models transfer Web knowledge to robotic control," in *Proc. of the Conf. on Robot Learning* (CoRL), 2023.
- [4] D. Ghosh, H. Walke, K. Pertsch, et al., "Octo: An open-source generalist robot policy," in *Proc. of Robotics: Science and Systems* (RSS), 2024.
- [5] K. Black, N. Brown, D. Driess, *et al.*, " π_0 : A vision-language-action flow model for general robot control," https://arxiv.org/abs/2410.24164, 2024.
- [6] D. Driess, F. Xia, M. S. M. Sajjadi, et al., "PaLM-e: An embodied multimodal language model," in Proc. of the Int. Conf. on Machine Learning (ICML), 2023, pp. 8469–8488.
- [7] M. J. Kim, K. Pertsch, S. Karamcheti, et al., "OpenVLA: An open-source vision-language-action model," in Proc. of the Conf. on Robot Learning (CoRL), 2024.
- [8] C.-L. Cheang, G. Chen, Y. Jing, et al., "GR-2: A generative videolanguage-action model with web-scale knowledge for robot manipulation," https://arxiv.org/abs/2410.06158, 2024.
- [9] A. Xie, L. Lee, T. Xiao, and C. Finn, "Decomposing the generalization gap in imitation learning for visual robotic manipulation," in *Proc. of* the IEEE Int. Conference on Robotics & Automation (ICRA), 2024.

- [10] C. Xu, Q. Li, J. Luo, and S. Levine, "RLDG: Robotic generalist policy distillation via reinforcement learning," https://arxiv.org/abs/ 2412.09858, 2024.
- [11] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *Proc. of the IEEE/RSJ Int. Conference on Intelligent Robots and Systems (IROS)*, 2012.
- [12] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox, "Gpu-accelerated robotic simulation for distributed reinforcement learning," in *Proc. of the Conf. on Robot Learning (CoRL)*, 2018.
- [13] F. Xiang, Y. Qin, K. Mo, et al., "SAPIEN: A simulated part-based interactive environment," in Proc. of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2020
- [14] M. Andrychowicz, A. Raichuk, P. Stańczyk, et al., "What matters for on-policy deep actor-critic methods? a large-scale study," in Proc. of the Int. Conf. on Learning Representations (ICLR), 2021.
- [15] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," https://arxiv.org/abs/1707. 06347, 2017.
- [16] A. O'Neill, A. Rehman, A. Gupta, et al., "Open x-embodiment: Robotic learning datasets and rt-x models," in Proc. of the IEEE Int. Conference on Robotics & Automation (ICRA), 2024.
- [17] R. Doshi, H. R. Walke, O. Mees, S. Dasari, and S. Levine, "Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation," in *Proc. of the Conf. on Robot Learning* (CoRL), 2024.
- [18] Y. Jiang, A. Gupta, Z. Zhang, et al., "VIMA: Robot manipulation with multimodal prompts," in Proc. of the Int. Conf. on Machine Learning (ICML), 2023.
- [19] W. M. Czarnecki, R. Pascanu, S. Osindero, S. Jayakumar, G. Swirszcz, and M. Jaderberg, "Distilling policy distillation," in *Proc. of the Int. Conf. on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- [20] A. A. Rusu, S. G. Colmenarejo, C. Gulcehre, et al., "Policy distillation," in Proc. of the Int. Conf. on Learning Representations (ICLR), 2016
- [21] S. Green, C. M. Vineyard, and Çetin Kaya Koç, "Distillation strategies for proximal policy optimization," https://arxiv.org/abs/1901.08128, 2019.
- [22] R. Agarwal, M. Schwarzer, P. S. Castro, A. C. Courville, and M. Bellemare, "Reincarnating reinforcement learning: Reusing prior computation to accelerate progress," in *Advances in Neural Information Processing Systems*, 2022.
- [23] G. Spigler, "Proximal policy distillation," https://arxiv.org/abs/2407. 15134, 2024.
- [24] D. Pomerleau, "ALVINN: An autonomous land vehicle in a neural network," in Advances in Neural Information Processing Systems, 1989.
- [25] A. Rajeswaran, V. Kumar, A. Gupta, et al., "Learning complex dexterous manipulation with deep reinforcement learning and demonstrations," in Proc. of Robotics: Science and Systems (RSS), 2018.
- [26] G. Libardi, G. De Fabritiis, and S. Dittert, "Guided exploration with proximal policy optimization using a single demonstration," in *Proc. of* the Int. Conf. on Machine Learning (ICML), 2021.
- [27] S. Levine, A. Kumar, G. Tucker, and J. Fu, "Offline reinforcement learning: Tutorial, review, and perspectives on open problems," https://arxiv.org/abs/2005.01643, 2020.
- [28] P. J. Ball, L. Smith, I. Kostrikov, and S. Levine, "Efficient online reinforcement learning with offline data," in *Proc. of the Int. Conf. on Machine Learning (ICML)*, 2023.
- [29] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Proc. of the Int. Conf. on Machine Learning (ICML)*, 2018.
- [30] R. S. Sutton and A. G. Barto, Reinforcement learning: an introduction. The MIT Press. 2018.
- [31] J. Gu, F. Xiang, X. Li, et al., "ManiSkill2: A unified benchmark for generalizable manipulation skills," in Proc. of the Int. Conf. on Learning Representations (ICLR), 2023.
- [32] S. Huang, R. F. J. Dossa, C. Ye, et al., "CleanRL: High-quality single-file implementations of deep reinforcement learning algorithms," Journal of Machine Learning Research, vol. 23, no. 274, pp. 1–18, 2022.
- [33] W. Zhao, J. P. Queralta, and T. Westerlund, "Sim-to-real transfer in deep reinforcement learning for robotics: a survey," in *Proc. of the IEEE Symposium Series on Computational Intelligence (SSCI)*, 2020.