

ON-POLICY POLICY GRADIENT REINFORCEMENT LEARNING WITHOUT ON-POLICY SAMPLING

Anonymous authors

Paper under double-blind review

ABSTRACT

On-policy reinforcement learning (RL) algorithms perform policy updates using i.i.d. trajectories collected by the current policy. However, after observing only a finite number of trajectories, on-policy sampling may produce data that fails to match the expected on-policy data distribution. This *sampling error* leads to noisy updates and data inefficient on-policy learning. Recent work in the policy evaluation setting has shown that non-i.i.d., off-policy sampling can produce data with lower sampling error than on-policy sampling can produce (Zhong et al., 2022). Motivated by this observation, we introduce an adaptive, off-policy sampling method to improve the data efficiency of on-policy policy gradient algorithms. Our method, **Proximal Robust On-Policy Sampling** (PROPS), reduces sampling error by collecting data with a *behavior policy* that increases the probability of sampling actions that are under-sampled with respect to the current policy. We empirically evaluate PROPS on both continuous-action MuJoCo benchmark tasks as well discrete-action tasks and demonstrate that (1) PROPS decreases sampling error throughout training and (2) improves the data efficiency of on-policy policy gradient algorithms.

1 INTRODUCTION

One of the most widely used classes of reinforcement learning (RL) algorithms is the class of on-policy policy gradient algorithms. These algorithms use gradient ascent on the parameters of a parameterized policy so as to increase the probability of observed actions with high expected returns under the current policy. The gradient is commonly estimated using the Monte Carlo estimator, an average computed over i.i.d. samples of trajectories from the current policy. The Monte Carlo estimator is consistent and unbiased; as the number of sampled trajectories increases, the empirical distribution of trajectories converges to the true distribution under the current policy, and thus the empirical gradient converges to the true gradient. However, the expense of environment interactions forces us to work with finite samples. Thus, the empirical distribution of the trajectories often differs from the desired on-policy data distribution. We refer to the mismatch between the empirical distribution of trajectories and the desired on-policy trajectory distribution as *sampling error*. This sampling error produces inaccurate gradient estimates, resulting in noisy policy updates, slower learning, and potentially convergence to suboptimal policies. With i.i.d. on-policy sampling, the only way to reduce sampling error is to collect more data.

Since on-policy sampling is so widely used to produce on-policy data, on-policy sampling is often taken to be an essential feature of data collection for on-policy learning (Silver, 2015; Achiam, 2018; Sutton and Barto, 2018). However, on-policy sampling is not explicitly required for on-policy learning; on-policy learning requires on-policy *data* – data whose state-conditioned empirical distribution of actions matches that of the current policy. On-policy sampling is a straightforward way to acquire on-policy data, though we can obtain such data more efficiently *without* on-policy sampling. To better illustrate this concept, consider an MDP with two discrete actions A and B, and suppose the current policy π places equal probability on both actions in some state s . When following π , after 10 visits to s , we may observe A 2 times and B 8 times rather than the expected 5 times. Alternatively, if we adaptively select the most under-sampled action upon every visit to s , we will observe each action an equal number of times. The first scenario illustrates on-policy sampling but not on-policy data; the second scenario uses *off-policy* sampling yet produces on-policy data.

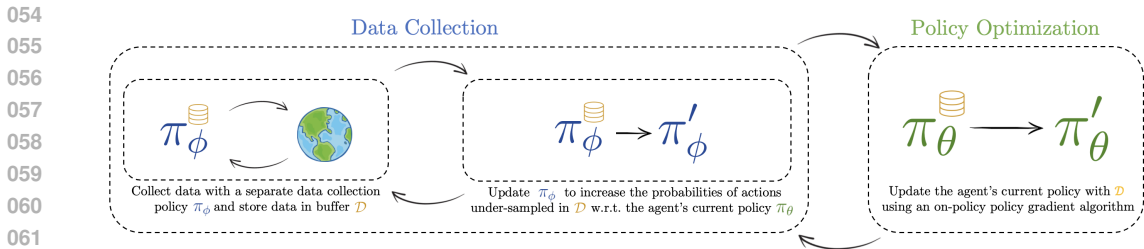


Figure 1: An overview of PROPS for on-policy policy gradient learning. Rather than collecting data \mathcal{D} via on-policy sampling from the agent’s current policy π_θ , we collect data with a separate data collection policy π_ϕ that we continually adapt to reduce sampling error in \mathcal{D} with respect to the agent’s current policy.

These observations raise the following question: can on-policy policy gradient algorithms learn more efficiently using on-policy data acquired *without* on-policy sampling? Recently, [Zhong et al. \(2022\)](#) showed that adaptive, off-policy sampling can yield data that more closely matches the on-policy distribution than data produced by i.i.d. on-policy sampling. However, this work was limited to the policy evaluation setting in which the on-policy distribution remains fixed. Turning from evaluation to control poses the challenge of a continually changing current policy.

In this work, we address this challenge and show for the first time that on-policy policy gradient algorithms are more data-efficient learners when they use on-policy data acquired with adaptive, off-policy sampling. Our method, **Proximal Robust On-Policy Sampling (PROPS)**¹, adaptively corrects sampling error in previously collected data by increasing the probability of sampling actions that are under-sampled with respect to the current policy. Fig. 1 provides an overview of PROPS. We empirically evaluate PROPS on continuous-action MuJoCo benchmark tasks as well as discrete action tasks and show that (1) PROPS reduces sampling error throughout training and (2) improves the data efficiency of on-policy policy gradient algorithms. In summary, our contributions are

1. We introduce an adaptive sampling algorithm that reduces sampling error in on-policy data collection.
2. We demonstrate empirically that our method improves the data efficiency of on-policy policy gradient algorithms and increases the fraction of training runs that converge to high-return policies.
3. Building off of the theoretical foundation laid by [Zhong et al. \(2022\)](#), this work improves the RL community’s understanding of a nuance in the on-policy vs off-policy dichotomy: on-policy learning requires on-policy data, not on-policy sampling.

2 RELATED WORK

Our work focuses on data collection in RL. In RL, data collection is often framed as an exploration problem, focusing on how an agent should explore its environment to efficiently learn an optimal policy. Prior RL works have proposed several exploration-promoting methods such as intrinsic motivation ([Pathak et al., 2017](#); [Sukhbaatar et al., 2018](#)), count-based exploration ([Tang et al., 2017](#); [Ostrovski et al., 2017](#)), and Thompson sampling ([Osband et al., 2013](#); [Sutton and Barto, 2018](#)). In contrast, our objective is to learn from the on-policy data distribution; we use adaptive data collection to more efficiently obtain this data distribution.

Prior works have used adaptive off-policy sampling to reduce sampling error in the policy evaluation subfield of RL. Most closely related is the work of [Zhong et al. \(2022\)](#) who first proposed that adaptive off-policy sampling could produce data that more closely matches the on-policy distribution than on-policy sampling could produce. [Mukherjee et al. \(2022\)](#) use a deterministic sampling rule to take actions in a particular proportion. Other bandit works use a non-adaptive exploration policy to collect additional data conditioned on previously collected data ([Tucker and Joachims, 2022](#); [Wan](#)

¹We include our codebase in the supplemental material.

et al., 2022; Konyushova et al., 2021). Since these works only focus on policy evaluation, they do not have to contend with a changing on-policy distribution as our work does for the control setting.

Several prior works propose importance sampling methods (Precup, 2000) to reduce sampling error without further data collection. In the RL setting, Hanna et al. (2021) showed that reweighting off-policy data according to an estimated behavior policy can correct sampling error and improve policy evaluation. Similar methods have been studied for temporal difference learning (Pavse et al., 2020) and policy evaluation in the bandit setting (Li et al., 2015; Narita et al., 2019). Conservative Data Sharing (Yu et al., 2021) reduces sampling error by selectively integrating offline data from multiple tasks. Our work instead focuses on using additional data collection to reduce sampling error.

As we will discuss in Section 5, the method we introduce permits data collected in one iteration of policy optimization to be re-used in future iterations rather than discarded as typically done by on-policy algorithms. Prior work has attempted to avoid discarding data by combining off-policy and on-policy updates with separate loss functions or by using alternative gradient estimates (Wang et al., 2016; Gu et al., 2016; 2017; Fakoor et al., 2020; O’Donoghue et al., 2016; Queeney et al., 2021). In contrast, our method modifies the sampling distribution at each iteration so that the entire data set of past and newly collected data matches the expected distribution under the current policy.

3 PRELIMINARIES

3.1 REINFORCEMENT LEARNING

We formalize the RL environment as a finite horizon Markov decision process (MDP) (Puterman, 2014) $(\mathcal{S}, \mathcal{A}, p, r, d_0, \gamma)$ with state space \mathcal{S} , action space \mathcal{A} , transition dynamics $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, initial state distribution d_0 , and reward discount factor $\gamma \in [0, 1)$. The state and action spaces may be discrete or continuous. We write $p(\cdot | s, a)$ to denote the distribution of next states after taking action a in state s . We consider stochastic policies $\pi_\theta : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ parameterized by θ , and we write $\pi_\theta(a|s)$ to denote the probability of sampling action a in state s and $\pi_\theta(\cdot|s)$ to denote the probability distribution over actions in state s . We additionally let $d_{\pi_\theta} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ denote the state-action visitation distribution, the distribution over state-action pairs induced by following π_θ . The RL objective is to find a policy that maximizes the expected sum of discounted rewards, defined as:

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

where the horizon H is the random variable representing the time-step when an episode ends. Throughout this paper, we refer to the policy used for data collection as the *behavior policy* and the policy trained to maximize its expected return as the *target policy*.

3.2 ON-POLICY POLICY GRADIENT ALGORITHMS

Policy gradient algorithms are one of the most widely used methods in RL. These methods perform gradient ascent over policy parameters to maximize an agent’s expected return $J(\theta)$ (Eq. 1). The gradient of the $J(\theta)$ with respect to θ , or *policy gradient*, is often given as:

$$\nabla_\theta J(\theta) = \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta} [A^{\pi_\theta}(s, a) \nabla_\theta \log \pi_\theta(a|s)], \quad (2)$$

where $A^{\pi_\theta}(s, a)$ is the *advantage* of choosing action a in state s and following π_θ thereafter. In practice, the expectation in Eq. 2 is approximated with Monte Carlo samples collected from π_θ and an estimate of A^{π_θ} used in place of the true advantages (Schulman et al., 2016). After updating the policy parameters with this estimated gradient, the previously collected trajectories \mathcal{D} become off-policy with respect to the updated policy. To ensure gradient estimation remains unbiased, on-policy algorithms discard historic data after each update and collect new data with the updated policy.

This foundational idea of policy learning via stochastic gradient ascent was first proposed by Williams (Williams, 1992) under the name REINFORCE. Since then, a large body of research has focused on developing more scalable policy gradient methods (Kakade, 2001; Schulman et al., 2015; Mnih et al., 2016; Espeholt et al., 2018; Lillicrap et al., 2015; Haarnoja et al., 2018). Arguably, the

most successful variant of policy gradient learning is proximal policy optimization (PPO) (Schulman et al., 2017), the algorithm of choice in several high-profile success stories (Berner et al., 2019; Akkaya et al., 2019; Vinyals et al., 2019). Rather than maximizing the standard RL objective (Eq. 1), PPO maximizes a surrogate objective:

$$\mathcal{L}_{\text{PPO}}(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}}) = \min(g(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}})A^{\pi_{\boldsymbol{\theta}_{\text{old}}}}(\mathbf{s}, \mathbf{a}), \text{clip}(g(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}}), 1 - \epsilon, 1 + \epsilon)A^{\pi_{\boldsymbol{\theta}_{\text{old}}}}(\mathbf{s}, \mathbf{a})), \quad (3)$$

where $\boldsymbol{\theta}_{\text{old}}$ denotes the policy parameters prior to the update, $g(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}})$ is the policy ratio $g(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}}) = \frac{\pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})}{\pi_{\boldsymbol{\theta}_{\text{old}}}(\mathbf{a}|\mathbf{s})}$, and the `clip` function with hyperparameter ϵ clips $g(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}})$ to the interval $[1 - \epsilon, 1 + \epsilon]$. This objective disincentivizes large changes to $\pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$. In contrast to other policy gradient algorithms which perform a single gradient update per data sample to avoid destructively large weight updates, PPO’s clipping mechanism allows the agent to perform multiple epochs of minibatch policy updates.

4 CORRECTING SAMPLING ERROR IN REINFORCEMENT LEARNING

In this section, we illustrate how sampling error can produce inaccurate policy gradient estimates and then describe how adaptive, off-policy sampling can reduce sampling error. For exposition, we assume finite state and action spaces. The policy gradient can then be written as:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \sum_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} d_{\pi_{\boldsymbol{\theta}}}^{\gamma}(\mathbf{s}, \mathbf{a}) [A^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})]. \quad (4)$$

The policy gradient is thus a linear combination of the gradient for each (\mathbf{s}, \mathbf{a}) pair $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$ weighted by $d_{\pi_{\boldsymbol{\theta}}}^{\gamma}(\mathbf{s}, \mathbf{a}) A^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}, \mathbf{a})$. Let \mathcal{D} be a dataset of trajectories. It is straightforward to show that the Monte Carlo estimate of the policy gradient can be written in a similar form as Equation 4 except with the true state-action visitation distribution replaced with the empirical visitation distribution, $d_{\mathcal{D}}(\mathbf{s}, \mathbf{a})$ (Hanna et al., 2021). Consequently, when (\mathbf{s}, \mathbf{a}) is over-sampled (*i.e.*, $d_{\mathcal{D}}(\mathbf{s}, \mathbf{a}) > d_{\pi_{\boldsymbol{\theta}}}^{\gamma}(\mathbf{s}, \mathbf{a})$), then $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$ contributes more to the overall gradient than it should. Similarly, when (\mathbf{s}, \mathbf{a}) is under-sampled, $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$ contributes less than it should.

We now provide a concrete example illustrating how small amounts of sampling error can cause the wrong actions to be reinforced, resulting in sub-optimal convergence. Suppose that in a particular state \mathbf{s}_0 , an agent places equal probability on two actions \mathbf{a}_0 and \mathbf{a}_1 with advantages $A^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}_0, \mathbf{a}_0) = 20$ and $A^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}_0, \mathbf{a}_1) = 15$, respectively. Since $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_0|\mathbf{s}_0) = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_1|\mathbf{s}_0)$ and $d_{\pi_{\boldsymbol{\theta}}}^{\gamma}(\mathbf{s}_0, \mathbf{a}_0) = d_{\pi_{\boldsymbol{\theta}}}^{\gamma}(\mathbf{s}_0, \mathbf{a}_1)$, the expected gradient will increase the probability of sampling the action with the larger advantage (\mathbf{a}_0). With on-policy sampling, after 10 visits to \mathbf{s}_0 , the agent will sample both actions 5 times in expectation. However, the agent may actually observe \mathbf{a}_0 4 times and \mathbf{a}_1 6 times. A Monte Carlo estimate of the policy gradient would then place $0.4 \cdot A^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}_0, \mathbf{a}_0) = 8$ weight on the gradient of \mathbf{a}_0 and $0.6 \cdot A^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}_0, \mathbf{a}_1) = 9$ weight on \mathbf{a}_1 , thus *decreasing* the probability of sampling the optimal \mathbf{a}_0 action.

Sampling error in on-policy sampling vanishes as the size of the batch of data used to estimate the gradient tends toward infinity. However, the preceding example suggests a simple strategy that would eliminate sampling error with finite data: have the agent adapt its probability on the next action it takes based on what actions it has already sampled. Continuing with our example, suppose the agent has visited \mathbf{s}_0 9 times and sampled \mathbf{a}_0 4 times and \mathbf{a}_1 5 times. With on-policy sampling, the agent may observe \mathbf{a}_1 again upon the next visit to \mathbf{s}_0 . Alternatively, the agent could sample its next action from a distribution that puts probability 1 on \mathbf{a}_0 and consequently produce an aggregate batch of data that contains both actions in their expected frequency. While this adaptive method is an off-policy sampling method, it produces data that exactly matches the on-policy distribution and will thus produce a more accurate gradient.

This example suggests that we can heuristically reduce sampling error by taking the most under-sampled action at a given state. Under a strong assumption that the MDP had a DAG structure, Zhong et al. (2022) proved that this heuristic results in the empirical distribution of states and actions in a fixed-horizon MDP converging to $d_{\pi_{\boldsymbol{\theta}}}(\mathbf{s}, \mathbf{a})$ and moreover converging at a faster rate than on-policy sampling. We remove this limiting assumption with the following result:

Proposition 1. Assume that data is collected with an adaptive behavior policy that always takes the most under-sampled action in each state, s , with respect to policy π , i.e. $a \leftarrow \arg \max_{a'} (\pi(a'|s) - \pi_{\mathcal{D}}(a'|s))$, where $\pi_{\mathcal{D}}$ is the empirical policy after m state-action pairs have been collected. Assume that \mathcal{S} and \mathcal{A} are finite and that the Markov chain induced by π is irreducible. Then we have that the empirical state visitation distribution, d_m , converges to the state distribution of π , d_π , with probability 1:

$$\forall s, \lim_{m \rightarrow \infty} d_m(s) = d_\pi(s).$$

Proof. See Appendix A. □

While adaptively sampling the most under-sampled action can reduce sampling error, this heuristic is difficult to implement in practice. In tasks with continuous states and actions, the $\arg \max$ often has no closed-form solution, and the empirical policy can be expensive to compute at every timestep. Building upon the concepts discussed in this section, the following section presents a *scalable* adaptive sampling algorithm that reduces sampling error in on-policy policy gradient learning.

5 PROXIMAL ROBUST ON-POLICY SAMPLING FOR POLICY GRADIENT ALGORITHMS

Our goal is to develop an adaptive, off-policy sampling algorithm that reduces sampling error in on-policy data collection for on-policy policy gradient algorithms. We outline a general framework for on-policy learning with an adaptive behavior policy in Algorithm 1. In this framework, the behavior policy π_ϕ and target policy π_θ are initially the same. The behavior policy collects a batch of m transitions, adds the batch to a data buffer \mathcal{D} , and then updates its weights such that the next batch it collects reduces sampling error in \mathcal{D} with respect to the target policy π_θ (Lines 7-10). Every n steps (with $n > m$), the agent updates its target policy with data from \mathcal{D} (Line 11). We refer to m and n as the *behavior batch size* and the *target batch size*, respectively.

A subtle implication of adaptive sampling is that it can correct sampling error in *any* empirical data distribution – even one generated by a different policy. Rather than discarding off-policy data from old policies – as is commonly done in on-policy learning – we let the data buffer hold up to b target batches (bn transitions) and call b the *buffer size*. If $b > 1$, then \mathcal{D} will contain historic off-policy data used in previous target policy updates. Regardless of how b is set, the role of the behavior policy is to continually adjust action probabilities for new samples so that the aggregate data distribution of \mathcal{D} matches the expected on-policy distribution of the current target policy (Line 10). Implementing Line 10 is the core challenge we address in the remainder of this section.

To ensure that the empirical distribution of \mathcal{D} matches the expected on-policy distribution, updates to π_ϕ should attempt to increase the probability of actions which are currently under-sampled with respect to π_θ . Zhong et al. (2022) recently developed a simple method called Robust On-policy Sampling (ROS) for making such updates. In particular, the gradient $\nabla_\phi \mathcal{L} := -\nabla_\phi \sum_{(s,a) \in \mathcal{D}} \log \pi_\phi(a|s)$ when evaluated at $\phi = \theta$ provides a direction to change ϕ such that under-sampled actions have their probabilities increased. Thus a single step of gradient ascent will increase the probability of under-sampled actions.² In theory and in simple RL policy evaluation

Algorithm 1 On-policy policy gradient algorithm with adaptive sampling

- 1: **Inputs:** Target batch size n , behavior batch size m , buffer size b .
 - 2: **Output:** Target policy parameters θ .
 - 3: Initialize target policy parameters θ .
 - 4: Initialize behavior policy parameters $\phi \leftarrow \theta$.
 - 5: Initialize empty buffer \mathcal{D} with capacity bn .
 - 6: **for** target update $i = 1, 2, \dots$ **do**
 - 7: **for** behavior update $j = 1, \dots, \lfloor n/m \rfloor$ **do**
 - 8: Collect batch of data \mathcal{B} by running π_ϕ .
 - 9: Append \mathcal{B} to buffer \mathcal{D} .
 - 10: Update π_ϕ with \mathcal{D} using Algorithm 2.
 - 11: Update π_θ with \mathcal{D} .
 - 12: **return** θ
-

²To add further intuition for this update, note that it is the opposite of a gradient ascent step on the log likelihood of \mathcal{D} . When starting at θ , gradient ascent on the data log likelihood will increase the probability of actions that are over-sampled relative to π_θ . Hence, the ROS update changes ϕ in the opposite direction.

tasks, this update was shown to improve the rate at which the empirical data distribution converges to the on-policy distribution – even when the empirical data distribution contains off-policy data. Unfortunately, there are two main challenges that render ROS unsuitable for Line 10 in Algorithm 1.

Challenge 1: Destructively large policy updates. Since the buffer \mathcal{D} may contain data collected from older target policies, some samples in \mathcal{D} may be very off-policy with respect to the current target policy such that $\log \pi_\phi(\mathbf{a}|\mathbf{s})$ is large and negative. Since $\nabla_\phi \log \pi_\phi(\mathbf{a}|\mathbf{s})$ increases in magnitude as $\pi_\phi(\mathbf{a}|\mathbf{s})$ tends towards zero, ROS incentivizes the agent to continually decrease the probability of these actions despite being extremely unlikely under the current target policy. Thus, off-policy samples can produce destructively large policy updates.

Challenge 2: Improper handling of continuous actions. In a continuous-action task, ROS may produce behavior policies that *increase* sampling error. A continuous-action task policy $\pi_\theta(\mathbf{a}|\mathbf{s})$ is typically parameterized as a Gaussian $\mathcal{N}(\boldsymbol{\mu}(\mathbf{s}), \Sigma(\mathbf{s}))$ with mean $\boldsymbol{\mu}(\mathbf{s})$ and diagonal covariance matrix $\Sigma(\mathbf{s})$. Since actions in the tail of the Gaussian far from the mean will always be under-sampled, the ROS update will continually push the components of $\boldsymbol{\mu}(\mathbf{s})$ towards $\pm\infty$ and the diagonal components of $\Sigma(\mathbf{s})$ towards 0 to increase the probability of sampling these actions. The result is a degenerate behavior policy that is so far from the target policy that sampling from it increases sampling error.³ We illustrate this scenario with 1-dimensional continuous actions in Fig. 6 of Appendix B.

To address these challenges, we propose a new behavior policy update. To address Challenge 1, first observe that the gradient of the ROS loss $\nabla_\phi \mathcal{L} = \nabla_\phi \log \pi_\phi(\mathbf{a}|\mathbf{s})|_{\phi=\theta}$ is equivalent to the policy gradient (Eq. 2) with $A^{\pi_\theta}(\mathbf{s}, \mathbf{a}) = -1, \forall(\mathbf{s}, \mathbf{a})$. Since the clipped surrogate objective of PPO (Eq. 3) prevents destructively large updates in on-policy policy gradient learning, we use a similar clipped surrogate objective in place of the ROS objective:

$$\mathcal{L}_{\text{CLIP}}(\mathbf{s}, \mathbf{a}, \phi, \theta, \epsilon_{\text{PROPS}}) = \min \left[-\frac{\pi_\phi(\mathbf{a}|\mathbf{s})}{\pi_\theta(\mathbf{a}|\mathbf{s})}, -\text{clip} \left(\frac{\pi_\phi(\mathbf{a}|\mathbf{s})}{\pi_\theta(\mathbf{a}|\mathbf{s})}, 1 - \epsilon_{\text{PROPS}}, 1 + \epsilon_{\text{PROPS}} \right) \right]. \quad (5)$$

Table 1 in Appendix B summarizes the behavior of $\mathcal{L}_{\text{CLIP}}$. Intuitively, this objective is equivalent to the PPO objective (Eq. 3) with $A(\mathbf{s}, \mathbf{a}) = -1, \forall(\mathbf{s}, \mathbf{a})$ and incentivizes the agent to decrease the probability of observed actions by at most a factor of $1 - \epsilon_{\text{PROPS}}$. Let $g(\mathbf{s}, \mathbf{a}, \phi, \theta) = \frac{\pi_\phi(\mathbf{a}|\mathbf{s})}{\pi_\theta(\mathbf{a}|\mathbf{s})}$. When $g(\mathbf{s}, \mathbf{a}, \phi, \theta) < 1 - \epsilon_{\text{PROPS}}$, this objective is clipped at $-(1 - \epsilon_{\text{PROPS}})$. The loss gradient $\nabla_\phi \mathcal{L}_{\text{CLIP}}$ becomes zero, and the (\mathbf{s}, \mathbf{a}) pair has no effect on the policy update. When $g(\mathbf{s}, \mathbf{a}, \phi, \theta) > 1 + \epsilon_{\text{PROPS}}$, clipping does not apply, and the gradient $\nabla_\phi \mathcal{L}_{\text{CLIP}}$ points in a direction that decreases the probability of $\pi_\phi(\mathbf{a}|\mathbf{s})$. As in the PPO update, this clipping mechanism avoids destructively large policy updates and permits us to perform many epochs of minibatch updates with the same batch of data.

To address the second challenge and prevent degenerate behavior policies, we introduce an auxiliary loss that incentivizes the agent to minimize the KL divergence between the behavior policy and target policy at states in the observed data. The full PROPS objective is then:

$$\mathcal{L}_{\text{PROPS}}(\mathbf{s}, \mathbf{a}, \phi, \theta, \epsilon_{\text{PROPS}}, \lambda) = \mathcal{L}_{\text{CLIP}}(\mathbf{s}, \mathbf{a}, \phi, \theta) - \lambda D_{\text{KL}}(\pi_\theta(\cdot|\mathbf{s})||\pi_\phi(\cdot|\mathbf{s})) \quad (6)$$

where λ is a regularization coefficient quantifying a trade-off between maximizing $\mathcal{L}_{\text{PROPS}}$ and minimizing D_{KL} . We provide full pseudocode for the PROPS update in Algorithm 2. Like ROS, we set the behavior policy parameters ϕ equal to the target policy parameters at the start of each behavior update, and then make a local adjustment to ϕ to increase the probabilities of under-sampled actions.

³This challenge is specific to continuous-action tasks and does not arise in discrete-action tasks.

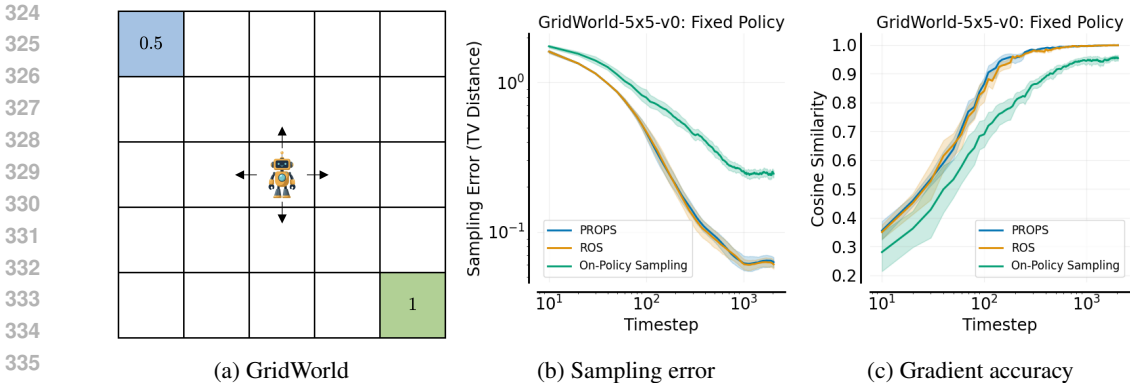


Figure 2: **(a)** A GridWorld task in which the agent receives reward +1 upon reaching the bottom right corner (the optimal goal), a reward of +0.5 upon reaching the top left corner (the suboptimal goal), and a reward of -0.01 . The agent always starts in the center of the grid. Under an initially uniform policy, the agent visits both goals with equal probability, and thus the true policy gradient increases the probability of reaching the optimal goal. **(b, c)** PROPS reduces sampling error and achieves more accurate gradients faster than on-policy sampling.

We stop the PROPS update when $D_{\text{KL}}(\pi_{\theta}||\pi_{\phi})$ reaches a chosen threshold δ_{PROPS} . This technique further safeguards against large policy updates and is used in widely adopted implementations of PPO (Raffin et al., 2021; Liang et al., 2018). The PROPS update allows us to efficiently learn a behavior policy that keeps the distribution of data in the buffer close to the expected distribution of the target policy.

6 EXPERIMENTS

The central goal of our work is to understand whether on-policy policy gradient algorithms are more data efficient learners when they use on-policy data acquired *without* on-policy sampling. Towards this goal, we design experiments to answer the two questions:

- Q1:** Does PROPS achieve lower sampling error than on-policy sampling during training?
Q2: Does PROPS increase the fraction of training runs that converge to high-return policies and improve the data efficiency of on-policy policy gradient algorithms?

Our empirical analysis focuses on continuous-state continuous-action MuJoCo benchmark tasks and a tabular 5x5 GridWorld task (Fig. 2a). We additionally consider three continuous-state discrete-action tasks: CartPole-v1, LunarLander-v2, and Discrete2D100-v0 – a 2D navigation task with 100 discrete actions. Due to space constraints, we include these tasks in Appendix D.4.

6.1 CORRECTING SAMPLING ERROR FOR A FIXED TARGET POLICY

We first study how quickly PROPS decreases sampling error when the target policy is fixed. This setting is similar to the policy evaluation setting considered by Zhong et al. (2022). As such, we provide two baselines for comparison: on-policy sampling and ROS.

Sampling error metrics. In GridWorld, we compute sampling error as the total variation (TV) distance between the empirical state-action visitation $d_{\mathcal{D}}(s, a)$ distribution – denoting the proportion of times (s, a) appears in buffer \mathcal{D} – and the true state-action visitation distribution under the agent’s policy: $\sum_{(s,a) \in \mathcal{D}} |d_{\mathcal{D}}(s, a) - d_{\pi_{\theta}}(s, a)|$. In continuous MuJoCo tasks where it is difficult to compute $d_{\mathcal{D}}(s, a)$, we follow Zhong et al. (2022) and measure sampling error using the KL-divergence $D_{\text{KL}}(\pi_{\mathcal{D}}||\pi_{\theta})$ between the empirical policy $\pi_{\mathcal{D}}$ and the target policy π_{θ} . We estimate $\pi_{\mathcal{D}}$ as the maximum likelihood estimate under data in the buffer via stochastic gradient ascent. Further details on how we compute $\pi_{\mathcal{D}}$ are in Appendix C.

Since it is straightforward to compute the true policy gradient in the GridWorld task, we additionally investigate how sampling error reduction affects gradient estimation by measuring the cosine simi-

378 larity between the empirical policy gradient $\nabla_{\theta} \hat{J}(\theta)$ and the true policy gradient. As the empirical
 379 gradient aligns more closely with the true gradient, the cosine similarity approaches 1.
 380

381 **Experimental setup.** In all tasks, we use a
 382 buffer with capacity of $\lfloor T/2 \rfloor$ samples, where
 383 T is the total number of samples collected by
 384 the agent. Thus, we expect sampling error to
 385 decrease over the first $\lfloor T/2 \rfloor$ samples and then
 386 remain roughly constant afterwards once the
 387 buffer is full. We use randomly initialized tar-
 388 get policies. Further experimental details such
 389 as hyperparameter tuning are described in Ap-
 390 pendix D.1.

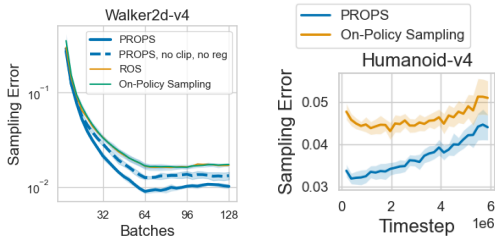
391 **Results.** As shown in Fig. 2b and 2c, in Grid-
 392 World, PROPS decreases sampling error faster
 393 than on-policy sampling, resulting in more ac-
 394 curate policy gradient estimates. PROPS and
 395 ROS to perform similarly, though this behavior
 396 is expected; in a tabular setting with a fixed tar-
 397 get policy (*i.e.* there is no off-policy data in the
 398 buffer), we do not encounter Challenge 1 and 2
 399 described in the previous section. In continu-
 400 ous MuJoCo tasks where Challenge 2 arises, PROPS
 401 decreases sampling error faster than on-policy
 402 sampling and ROS (Fig. 3a). In fact, ROS shows
 403 little to no improvement over on-policy sampling
 404 in every MuJoCo task. This limitation of ROS is
 405 unsurprising, as Zhong et al. (2022) showed that
 406 ROS struggled to reduce sampling error even in
 407 low-dimensional continuous-action tasks. More-
 408 over, PROPS decreases sampling error without
 409 clipping and regularization, emphasizing how
 410 adaptive off-policy sampling alone decreases
 411 sampling error. Due to space constraints, we in-
 412 clude results for the remaining environments in
 413 Appendix D.1. We additionally include experi-
 414 ments using a fixed, randomly initialized target
 415 policy as well as ablation studies isolating the
 416 effects of PROPS’s objective clipping and regu-
 417 larization in Appendix D.1. Results with a ran-
 418 dom target policy are qualitatively similar to
 419 those in Fig. 3a, and we observe that clipping
 420 and regularization both individually help re-
 421 duce sampling error.

412 6.2 CORRECTING SAMPLING ERROR DURING RL TRAINING

413 We are ultimately interested in understanding how
 414 replacing on-policy sampling with PROPS affects
 415 the data efficiency of on-policy learning, where
 416 the target policy is continually updated. In the
 417 following experiments, we train RL agents with
 418 PROPS and on-policy sampling to evaluate (1) the
 419 data efficiency of training, (2) the distribution
 420 of returns achieved at the end of training, and
 421 (3) the sampling error throughout training. We
 422 use the same sampling error metrics described in
 423 the previous section and measure data efficiency
 424 as the return achieved within a fixed training
 425 budget. Since ROS (Zhong et al., 2022) is
 426 computationally expensive and fails to reduce
 427 sampling error in MuJoCo tasks even with a
 428 fixed policy, we omit it from MuJoCo experi-
 429 ments.

429 **Experimental setup.** We use PPO (Schulman et al., 2017) to update the target policy. We consider
 430 two baseline methods for providing data to compute PPO updates: (1) vanilla PPO with on-policy
 431 sampling, and (2) PPO with on-policy sampling and a buffer of size b (PPO-BUFFER). PPO-BUFFER
 is a naive method for improving data efficiency of on-policy algorithms by reusing off-policy data
 collected by old target policies as if it were on-policy data. Although PPO-BUFFER computes biased
 gradients, it has been successfully applied in difficult learning tasks (Berner et al., 2019). Since
 PROPS and PPO-BUFFER have access to the same amount of data for each policy update, any per-
 formance difference between these two methods can be attributed to differences in how they sample
 actions during data collection.

In MuJoCo experiments, we set $b = 2$ such that agents retain a batch of data for one extra iteration
 before discarding it. In GridWorld, we use $b = 1$ and discard all historic data. Since PROPS and
 PPO-BUFFER compute target policy updates with b times as much learning data as PPO, we integrate
 this extra data by increasing the minibatch size for target and behavior policy updates by a factor



(a) Fixed target policy. (b) During RL training.

Figure 3: PROPS reduces sampling error faster than on-policy sampling and ROS. In (a), the ROS and on-policy sampling curves overlap. Solid curves denote means over 5 seeds. Shaded regions denote 95% confidence intervals.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

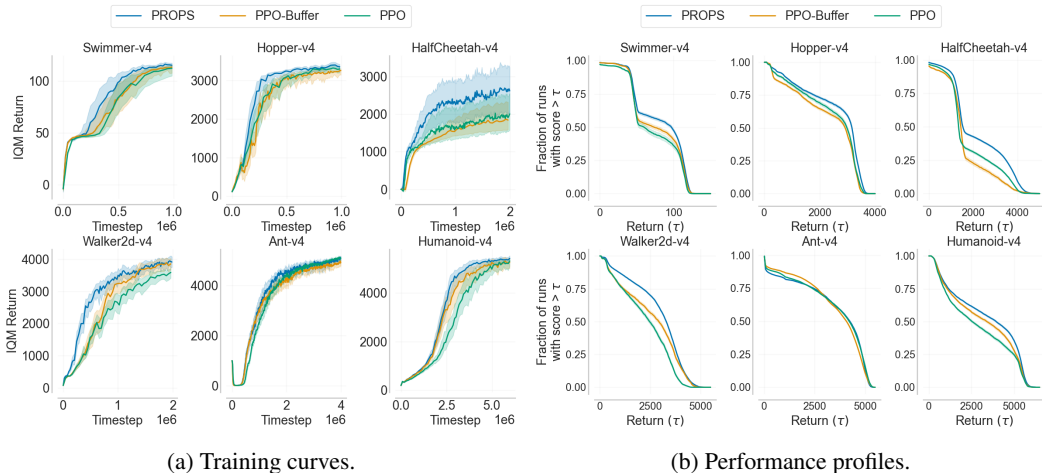
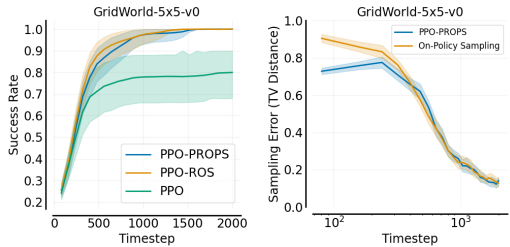


Figure 4: (a) IQM returns over 50 seeds. Shaded regions denote 95% bootstrap confidence intervals. (b) Performance profiles over 50 seeds. Higher values correspond to more reliable convergence to high-return policies. Shaded regions denote 95% bootstrap confidence intervals.

of *b*. Further experimental details including hyperparameter tuning are described in Appendix E. For MuJoCo tasks, we plot the interquartile mean (IQM) return throughout training as well as the distribution of returns achieved at the end of training (*i.e.*, the performance profile) (Agarwal et al., 2021). For GridWorld, we plot the agent’s success rate, the fraction of times it finds the optimal goal.

Results. As shown in Fig. 5a, on-policy sampling has approximately a 77% success rate on GridWorld, whereas PROPS and ROS achieve 100% success rate. In Fig. 4a, PROPS achieves higher return than both PPO and PPO-BUFFER throughout training in all MuJoCo tasks except Ant-v4 where PROPS dips slightly below PPO’s return near the end of training. Moreover, in Fig. 4b, the performance profile of PROPS almost always lies above the performance profiles of PPO and PPO-BUFFER, indicating that any given run of PROPS is more like to obtain a higher return than PPO-BUFFER. Thus, we affirmatively answer Q2 posed at the start of this section: PROPS increases the fraction of training runs with high return and increases data efficiency.

In Appendix D.4, we provide additional experiments demonstrating that PROPS improves data efficiency in discrete-action tasks. We additionally ablate the buffer size *b* in Appendix D.2. We find that data efficiency may decrease with a larger buffer size. Intuitively, the more historic data kept around, the more data that must be collected to impact the aggregate data distribution.



Having established that PROPS improves data efficiency, we now investigate if PROPS is appropriately adjusting the data distribution of the buffer by comparing the sampling error achieved throughout training with PROPS and PPO-BUFFER. Training with PROPS produces a different sequence of target policies than training with PPO-BUFFER produces. To provide a fair comparison, we compute sampling error for PPO-BUFFER using the target policy sequence produced by PROPS. More concretely, we fill a second buffer with on-policy samples collected by the *target policies* produced while training with PROPS and then compute the sampling error using data in this buffer.

As shown in Fig. 5b, PROPS achieves lower sampling error than on-policy sampling with a buffer in Humanoid-v4. Due to space constraints, we provide sampling error curves for the remaining MuJoCo environments in Appendix D.2. In GridWorld, PROPS and ROS reduce sampling error in

(a) Success rate. (b) Sampling error. Figure 5: GridWorld RL experiments over 50 seeds.

the first 300 steps and closely matches on-policy sampling afterwards. We use a batch size of 80 in these experiments, and as the target policy becomes more deterministic, larger batch sizes are needed to observe differences between PROPS and on-policy sampling.⁴

We additionally ablate the effects of the clipping coefficient ϵ_{PROPS} and regularization coefficient λ in Appendix D.2. Without clipping or without regularization, PROPS often achieves greater sampling error than on-policy sampling, indicating that both help to keep sampling error low. Moreover, data efficiency generally decreases when we remove clipping or regularization, showing both are essential to PROPS. Thus, we affirmatively answer **Q1** posed at the start of this section: PROPS achieves lower sampling error than on-policy sampling when the target policy is fixed and during RL training.

7 DISCUSSION

This work has shown that adaptive, off-policy sampling can be used to reduce sampling error in data collected throughout RL training and improve the data efficiency of on-policy policy gradient algorithms. We have introduced an algorithm that scales adaptive off-policy sampling to continuous control RL benchmarks and enables tracking of the changing on-policy distribution. By integrating this data collection procedure into the popular PPO algorithm, the main conclusion of our analysis is that on-policy learning algorithms learn most efficiently with on-policy data, *not* on-policy sampling. In this section, we discuss limitations of our work and present opportunities for future research.

PROPS builds upon the ROS algorithm of [Zhong et al. \(2022\)](#). While [Zhong et al. \(2022\)](#) focused on theoretical analysis and policy evaluation in small scale domains, we chose to focus on empirical analysis with policy learning in standard RL benchmarks. An important direction for future work would be theoretical analysis of PROPS, in particular whether PROPS also enjoys the same faster convergence rate that was shown for ROS relative to on-policy sampling.

A limitation of PROPS is that the update indiscriminately increases the probability of under-sampled actions without considering their importance in gradient computation. For instance, if an under-sampled action has zero advantage, it has no impact on the gradient and need not be sampled. An interesting direction for future work could be to prioritize correcting sampling error for (s, a) that have the largest influence on the gradient estimate, *i.e.*, large advantage (positive or negative).

Beyond these more immediate directions, our work opens up other opportunities for future research. A less obvious feature of the PROPS behavior policy update is that it can be used track the empirical data distribution of *any* desired policy, not only that of the current policy. This feature means PROPS has the potential to be integrated into off-policy RL algorithms and used so that the empirical distribution more closely matches a desired exploration distribution.^f Thus, PROPS could be used to perform focused exploration without explicitly tracking state and action counts.

8 CONCLUSION

In this paper, we ask whether on-policy policy gradient methods are more data efficient using on-policy sampling or on-policy *data* acquired *without* on-policy sampling. To answer this question, we introduce an adaptive, *off-policy* sampling method for on-policy policy gradient learning that collects data such that the empirical distribution of sampled actions closely matches the expected on-policy data distribution at observed states. Our method, Proximal Robust On-policy Sampling (PROPS), periodically updates the data collecting behavior policy so as to increase the probability of sampling actions that are currently under-sampled with respect to the on-policy distribution. Furthermore, rather than discarding collected data after every policy update, PROPS permits more data efficient on-policy learning by using data collection to adjust the distribution of previously collected data to be approximately on-policy. We replace on-policy sampling with PROPS to generate data for the widely-used PPO algorithm and empirically demonstrate that PROPS produces data that more closely matches the expected on-policy distribution and yields more data efficient learning compared to on-policy sampling.

⁴When the target policy is deterministic, we always have zero sampling error, and PROPS will exactly match on-policy sampling.

REFERENCES

- 540
541
542 Joshua Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- 543
544 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G Bellemare.
545 Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Informa-*
546 *tion Processing Systems*, 2021.
- 547
548 Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron,
549 Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a
robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- 550
551 Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy
552 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large
scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 553
554 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
555 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 556
557 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam
558 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with im-
559 portance weighted actor-learner architectures. In *International conference on machine learning*,
pages 1407–1416. PMLR, 2018.
- 560
561 Rasool Fakoor, Pratik Chaudhari, and Alexander J Smola. P3o: Policy-on policy-off policy opti-
562 mization. In *Uncertainty in Artificial Intelligence*, pages 1017–1027. PMLR, 2020.
- 563
564 Shixiang Gu, Timothy Lillicrap, Zoubin Ghahramani, Richard E Turner, and Sergey Levine. Q-prop:
565 Sample-efficient policy gradient with an off-policy critic. *arXiv preprint arXiv:1611.02247*, 2016.
- 566
567 Shixiang Shane Gu, Timothy Lillicrap, Richard E Turner, Zoubin Ghahramani, Bernhard Schölkopf,
568 and Sergey Levine. Interpolated policy gradient: Merging on-policy and off-policy gradient esti-
mation for deep reinforcement learning. *Advances in neural information processing systems*, 30,
2017.
- 569
570 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
571 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and appli-
cations. *arXiv preprint arXiv:1812.05905*, 2018.
- 572
573 Josiah P Hanna, Scott Niekum, and Peter Stone. Importance sampling in reinforcement learning
574 with an estimated behavior policy. *Machine Learning*, 110(6):1267–1317, 2021.
- 575
576 Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Ki-
577 nal Mehta, and João G.M. Araújo. Cleanrl: High-quality single-file implementations of deep
578 reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022.
URL <http://jmlr.org/papers/v23/21-1342.html>.
- 579
580 Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*,
14, 2001.
- 581
582 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua
583 Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations*,
584 *ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL
585 <http://arxiv.org/abs/1412.6980>.
- 586
587 Ksenia Konyushova, Yutian Chen, Thomas Paine, Caglar Gulcehre, Cosmin Paduraru, Daniel J
588 Mankowitz, Misha Denil, and Nando de Freitas. Active offline policy selection. *Advances in*
Neural Information Processing Systems, 34:24631–24644, 2021.
- 589
590 Lihong Li, Rémi Munos, and Csaba Szepesvári. Toward minimax off-policy value estimation. In
591 *Artificial Intelligence and Statistics*, pages 608–616. PMLR, 2015.
- 592
593 Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gon-
zalez, Michael Jordan, and Ion Stoica. Rllib: Abstractions for distributed reinforcement learning.
In *International Conference on Machine Learning*, pages 3053–3062. PMLR, 2018.

- 594 Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
595 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv*
596 *preprint arXiv:1509.02971*, 2015.
597
- 598 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim
599 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
600 learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
601
- 602 Subhojyoti Mukherjee, Josiah P Hanna, and Robert D Nowak. Revar: Strengthening policy evalu-
603 ation via reduced variance sampling. In *Uncertainty in Artificial Intelligence*, pages 1413–1422.
604 PMLR, 2022.
- 605 Yusuke Narita, Shota Yasui, and Kohei Yata. Efficient counterfactual learning from bandit feedback.
606 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4634–4641,
607 2019.
608
- 609 Brendan O’Donoghue, Rémi Munos, Koray Kavukcuoglu, and Volodymyr Mnih. Pqg: Combining
610 policy gradient and q-learning. *ArXiv*, abs/1611.01626, 2016.
611
- 612 Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via
613 posterior sampling. *Advances in Neural Information Processing Systems*, 26, 2013.
- 614 Georg Ostrovski, Marc G Bellemare, Aäron Oord, and Rémi Munos. Count-based exploration
615 with neural density models. In *International conference on machine learning*, pages 2721–2730.
616 PMLR, 2017.
617
- 618 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by
619 self-supervised prediction. In *International conference on machine learning*, pages 2778–2787.
620 PMLR, 2017.
- 621 Brahma S. Pavse, Ishan Durugkar, Josiah P. Hanna, and Peter Stone. Reducing sampling error in
622 batch temporal difference learning. In *International Conference on Machine Learning*, 2020.
623
- 624 Doina Precup. Eligibility traces for off-policy policy evaluation. *Computer Science Department*
625 *Faculty Publication Series*, page 80, 2000.
626
- 627 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
628 Wiley & Sons, 2014.
629
- 630 James Queeney, Ioannis Ch. Paschalidis, and Christos G. Cassandras. Generalized proximal policy
631 optimization with sample reuse. In *Advances in Neural Information Processing Systems*, vol-
632 ume 34. Curran Associates, Inc., 2021.
- 633 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
634 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
635 *Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/20-1364.html)
636 [20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
637
- 638 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
639 policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR,
640 2015.
- 641 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-
642 dimensional continuous control using generalized advantage estimation. In *International Con-*
643 *ference on Learning Representations (ICLR)*, 2016.
644
- 645 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
646 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
647
- David Silver. RL course by david silver, lecture 5: Model-free control. 2015.

- 648 Sainbayar Sukhbaatar, Zeming Lin, Ilya Kostrikov, Gabriel Synnaeve, Arthur Szlam, and Rob
649 Fergus. Intrinsic motivation and automatic curricula via asymmetric self-play. In *6th Inter-*
650 *national Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April*
651 *30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL [https://](https://openreview.net/forum?id=SKT5Yg-RZ)
652 openreview.net/forum?id=SKT5Yg-RZ.
- 653 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 654
- 655 Haoran Tang, Rein Houthoofd, Davis Foote, Adam Stooke, OpenAI Xi Chen, Yan Duan, John Schul-
- 656 man, Filip DeTurck, and Pieter Abbeel. # exploration: A study of count-based exploration for
- 657 deep reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- 658
- 659 Aaron David Tucker and Thorsten Joachims. Variance-optimal augmentation logging for counter-
- 660 factual evaluation in contextual bandits. *arXiv preprint arXiv:2202.01721*, 2022.
- 661
- 662 Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Juny-
- 663 oung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster
- 664 level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 665
- 666 Runzhe Wan, Branislav Kveton, and Rui Song. Safe exploration for efficient policy evaluation and
- 667 comparison. In *International Conference on Machine Learning*, pages 22491–22511. PMLR, 2022.
- 668
- 669 Ziyu Wang, Victor Bapst, Nicolas Heess, Volodymyr Mnih, Remi Munos, Koray Kavukcuoglu,
- 670 and Nando de Freitas. Sample efficient actor-critic with experience replay. *arXiv preprint*
arXiv:1611.01224, 2016.
- 671
- 672 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
- 673 learning. *Reinforcement learning*, pages 5–32, 1992.
- 674
- 675 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, and Chelsea Finn.
- 676 Conservative data sharing for multi-task offline reinforcement learning. *Advances in Neural In-*
formation Processing Systems, 34:11501–11516, 2021.
- 677
- 678 Rujie Zhong, Duohan Zhang, Lukas Schäfer, Stefano Albrecht, and Josiah Hanna. Robust on-
- 679 policy sampling for data-efficient policy evaluation in reinforcement learning. *Advances in Neural*
Information Processing Systems, 35:37376–37388, 2022.
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

Appendix

Table of Contents

A Theoretical Results	14
B PROPS Implementation Details	16
C Computing Sampling Error	17
D Additional Experiments	17
D.1 Correcting Sampling Error for a Fixed Target Policy	18
D.2 Correcting Sampling Error During RL Training	19
D.3 Bias and Variance of PROPS	20
D.4 Discrete-Action Tasks	21
D.5 Runtime Comparisons	21
E Hyperparameter Tuning for RL Training	22

A THEORETICAL RESULTS

In this section, we present the proof of Proposition 2. We use d_m , π_m , and p_m as the empirical state visitation distribution, empirical policy, and empirical transition probabilities after m state-action pairs have been taken, respectively. That is, $d_m(s)$ is the proportion of the m states that are s , $\pi_m(a|s)$ is the proportion of the time that action a was observed in state s , and $p_m(s'|s, a)$ is the proportion of the time that the state changed to s' after action a was taken in state s .

Proposition 2. *Assume that data is collected with an adaptive behavior policy that always takes the most under-sampled action in each state, s , with respect to policy π , i.e., $a \leftarrow \arg \max_{a'} (\pi(a'|s) - \pi_m(a'|s))$. We further assume that \mathcal{S} and \mathcal{A} are finite. Then we have that the empirical state visitation distribution, d_m , converges to the state distribution of π , d_π , with probability 1:*

$$\forall s, \lim_{m \rightarrow \infty} d_m(s) = d_\pi(s).$$

Proof. The proof of this theorem builds upon Lemma 1 and 2 by Zhong et al. (2022). Note that these lemmas superficially concern the ROS method whereas we are interested in data collection by taking the most under-sampled action at each step. However, as stated in the proof by Zhong et al. (2022), these methods are equivalent under an assumption they make about the step-size parameter of the ROS method. Thus, we can immediately adopt these lemmas for this proof.

Under Lemma 1 of Zhong et al. (2022), we have that $\lim_{m \rightarrow \infty} \pi_m(a|s) = \pi(a|s)$ for any state s under this adaptive data collection procedure. We then have the following $\forall s$:

$$\begin{aligned} \lim_{m \rightarrow \infty} d_m(s) &\stackrel{(a)}{=} \lim_{m \rightarrow \infty} \sum_{\tilde{s}} \sum_{\tilde{a}} p_m(s|\tilde{s}, \tilde{a}) \pi_m(\tilde{a}|\tilde{s}) d_m(\tilde{s}) \\ &= \sum_{\tilde{s}} \sum_{\tilde{a}} \lim_{m \rightarrow \infty} p_m(s|\tilde{s}, \tilde{a}) \pi_m(\tilde{a}|\tilde{s}) d_m(\tilde{s}) \\ &= \sum_{\tilde{s}} \sum_{\tilde{a}} \lim_{m \rightarrow \infty} p_m(s|\tilde{s}, \tilde{a}) \lim_{m \rightarrow \infty} \pi_m(\tilde{a}|\tilde{s}) \lim_{m \rightarrow \infty} d_m(\tilde{s}) \\ &\stackrel{(b)}{=} \sum_{\tilde{s}} \sum_{\tilde{a}} p(s|\tilde{s}, \tilde{a}) \pi(\tilde{a}|\tilde{s}) \lim_{m \rightarrow \infty} d_m(\tilde{s}). \end{aligned}$$

Here, (a) follows from the fact that the empirical frequency of state s can be obtained by considering all possible transitions that lead to s . The last line, (b), holds with probability 1 by the strong law of large numbers and Lemma 2 of [Zhong et al. \(2022\)](#).

We now have a system of $|\mathcal{S}|$ variables and $|\mathcal{S}|$ linear equations. Define variables $x(s) := \lim_{m \rightarrow \infty} d_m(s)$ and let $\mathbf{x} \in \mathbf{R}^{|\mathcal{S}|}$ be the vector of these variables. We then have $x = P^\pi x$ where $P^\pi \in \mathbf{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is the transition matrix of the Markov chain induced by running policy π . Assuming that this Markov chain is irreducible, d_π is the unique solution to this system of equations and hence $\lim_{m \rightarrow \infty} d_m(s) = d_\pi(s), \forall s$.

□

Next, we provide additional theory to describe the relationship between different hyperparameters in PROPS:

1. The amount of sampling error in previously collected data and the size of behavior policy updates.
2. The amount of historic data retained by an agent and the amount of additional data the behavior policy must collect to reduce sampling error.

For simplicity, we first focus on a simple bandit setting and then extend to a tabular RL setting.

Suppose we have already collected m state-action pairs and these have been observed with empirical distribution $\pi_m(\mathbf{a})$. From what distribution should we sample an additional k state-action pairs so that the empirical distribution over the $m + k$ samples is equal in expectation to π_θ ?

Proposition 3. *Assume that m actions have been collected by running some policy $\pi_\theta(\mathbf{a})$ and $\pi_m(\mathbf{a})$ is the empirical distribution on this dataset. If we collect an additional k state-action pairs using the following distribution, and if $(m + k)\pi_\theta(\mathbf{a}) \geq m \cdot \pi_m(\mathbf{a})$, then the aggregate empirical distribution over the $m + k$ pairs is equal to $\pi_\theta(\mathbf{a})$ in expectation:*

$$\pi_b(\mathbf{a}) := \frac{1}{Z} \left[\pi_\theta(\mathbf{a}) + \frac{m}{k} (\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})) \right]$$

where $Z = \sum_{\mathbf{a} \in \mathcal{A}} \left[\pi_\theta(\mathbf{a}) + \frac{m}{k} (\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})) \right]$ is a normalization coefficient.

Proof. Observe that $(m + k)\pi_\theta(\mathbf{a})$ is the expected number of times \mathbf{a} is sampled under π_θ after $m + k$ steps, $m \cdot \pi_m(\mathbf{a})$ is the number of times each \mathbf{a} was sampled thus far, and $k \cdot \pi_b(\mathbf{a})$ is the expected number of times \mathbf{a} is sampled under our behavior policy after k steps. We want to choose $\pi_b(\mathbf{a})$ such that $(m + k)\pi_\theta(\mathbf{a}) = m \cdot \pi_m(\mathbf{a}) + k \cdot \pi_b(\mathbf{a})$ in expectation.

$$\begin{aligned} (m + k)\pi_\theta(\mathbf{a}) &= k \cdot \pi_b(\mathbf{a}) + m \cdot \pi_m(\mathbf{a}) \\ -k \cdot \pi_b(\mathbf{a}) &= m \cdot \pi_m(\mathbf{a}) - (m + k)\pi_\theta(\mathbf{a}) \\ \pi_b(\mathbf{a}) &= -\frac{m}{k} \pi_m(\mathbf{a}) + \left(\frac{m}{k} + 1 \right) \pi_\theta(\mathbf{a}) \\ &= \pi_\theta(\mathbf{a}) + \frac{m}{k} (\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})) \end{aligned}$$

Note that $\pi_b(\mathbf{a})$ will be a valid probability distribution after normalizing only if

$$\begin{aligned} \pi_\theta(\mathbf{a}) + \frac{m}{k} (\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})) &\geq 0 \\ \left(\frac{m}{k} + 1 \right) \pi_\theta(\mathbf{a}) &\geq \frac{m}{k} \pi_m(\mathbf{a}) \\ (m + k) \pi_\theta(\mathbf{a}) &\geq m \cdot \pi_m(\mathbf{a}). \end{aligned}$$

If $(m + k) \pi_\theta(\mathbf{a}) < m \cdot \pi_m(\mathbf{a})$, then prior to collecting additional data with our behavior policy, \mathbf{a} already appears in our data more times in our data than it would in expectation after $m + k$ steps under π_θ . In other words, we would need to collect more than k additional samples to achieve zero sampling error (or discard some previously collected samples).

□

810 **When sampling error is large, behavior policy updates must also be large.** Intuitively, the
 811 difference $\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})$ is the mismatch between the true and empirical visitation distributions,
 812 so adding this term to d_{π_θ} adjusts d_{π_θ} to reduce this mismatch. If $\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a}) < 0$, then \mathbf{a} is over-
 813 sampled w.r.t π_θ , and π_b will decrease the probability of sampling \mathbf{a} . If $\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a}) > 0$, then \mathbf{a} is
 814 under-sampled w.r.t π_θ , and π_b will increase the probability of sampling \mathbf{a} . When $|\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})|$
 815 is small, the optimal $\pi_b(\mathbf{a})$ requires only a small adjustment from π_θ (i.e., a small update to the
 816 behavior policy is sufficient to reduce sampling error). When $|\pi_\theta(\mathbf{a}) - \pi_m(\mathbf{a})|$ is large, the optimal
 817 $\pi_b(\mathbf{a})$ requires a large adjustment from π_θ (i.e., a large to the behavior policy is needed to reduce
 818 sampling error).

819 **When we retain large amounts of historic data, the behavior policy must collect a large amount**
 820 **of additional data to reduce sampling error in the aggregate distribution.** The $\frac{m}{k}$ factor implies
 821 that how much we adjust d_{π_θ} depends on how much data we have already collected (m) and how
 822 much additional data we will collect (k). If the k additional samples to collect represent a small
 823 fraction of the aggregate $m + k$ samples (i.e. $k \ll m$), then $\frac{m}{k}$ is large, and the adjustment to d_{π_θ}
 824 is large. This case generally arises when we retain more and more historic data. If the k additional
 825 samples to collect represent a large fraction of the aggregate $m + k$ samples (i.e. $k \gg m$), then $\frac{m}{k}$
 826 is small, and the adjustment to d_{π_θ} is small. This case generally arises when we retain little to no
 827 historic data.

828 The next proposition extends this analysis to the tabular RL setting.

829 **Proposition 4.** *Assume that m state-action pairs have been collected by running some policy and*
 830 *$d_m(\mathbf{s}, \mathbf{a})$ is the empirical distribution on this dataset. If we collect an additional k state-action*
 831 *pairs using the following distribution, and if $(m + k)d_{\pi_\theta}(\mathbf{s}, \mathbf{a}) \geq m \cdot d_m(\mathbf{s}, \mathbf{a})$, then the aggregate*
 832 *empirical distribution over the $m + k$ pairs is equal to $d_{\pi_\theta}(\mathbf{s}, \mathbf{a})$ in expectation:*

$$833 \quad d_b(\mathbf{s}, \mathbf{a}) := \frac{1}{Z} \left[d_{\pi_\theta}(\mathbf{s}, \mathbf{a}) + \frac{m}{k} (d_{\pi_\theta}(\mathbf{s}, \mathbf{a}) - d_m(\mathbf{s}, \mathbf{a})) \right]$$

834 where $Z = \sum_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} \left[d_{\pi_\theta}(\mathbf{s}, \mathbf{a}) + \frac{m}{k} (d_{\pi_\theta}(\mathbf{s}, \mathbf{a}) - d_m(\mathbf{s}, \mathbf{a})) \right]$ is a normalization coefficient.

835 *Proof.* The proof is identical to the proof of Proposition 3, replacing $\pi_\theta(\mathbf{a})$, $\pi_m(\mathbf{a})$, and $\pi_b(\mathbf{a})$ with
 836 $d_{\pi_\theta}(\mathbf{s}, \mathbf{a})$, $d_m(\mathbf{s}, \mathbf{a})$, and $d_b(\mathbf{s}, \mathbf{a})$. \square

837 In practice, we cannot sample directly from the visitation distribution $d_b(\mathbf{s}, \mathbf{a})$ in Proposition 4
 838 and instead approximate sampling from this distribution by sampling from its corresponding policy
 839 $\pi_b(\mathbf{a}|\mathbf{s}) = d_b(\mathbf{s}, \mathbf{a}) / \sum_{(\mathbf{s}', \mathbf{a}') \in \mathcal{S} \times \mathcal{A}} d_b(\mathbf{s}', \mathbf{a}')$.

840 B PROPS IMPLEMENTATION DETAILS

841 In this appendix, we describe two relevant implementation details for the PROPS update (Algo-
 842 rithm 2). We additionally summarize the behavior of PROPS’s clipping mechanism in Table 1.

- 843 1. **PROPS update:** The PROPS update adapts the behavior policy to reduce sampling error in
 844 the buffer \mathcal{D} . When performing this update with a full buffer, we exclude the oldest batch
 845 of data collected by the behavior policy (i.e., the m oldest transitions in \mathcal{D}); this data will be
 846 evicted from the buffer before the next behavior policy update and thus does not contribute
 847 to sampling error in \mathcal{D} .
- 848 2. **Behavior policy class:** We compute behavior policies from the same policy class used for
 849 target policies. In particular, we consider Gaussian policies which output a mean $\mu(\mathbf{s})$ and
 850 a variance $\sigma^2(\mathbf{s})$ and then sample actions $\mathbf{a} \sim \pi(\cdot|\mathbf{s}) \equiv \mathcal{N}(\mu(\mathbf{s}), \sigma^2(\mathbf{s}))$. In principle, the
 851 target and behavior policy classes can be different. However, using the same class for both
 852 policies allows us to easily initialize the behavior policy equal to the target policy at the
 853 start of each update. This initialization is necessary to ensure the PROPS update increases
 854 the probability of sampling actions that are currently under-sampled with respect to the
 855 target policy.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

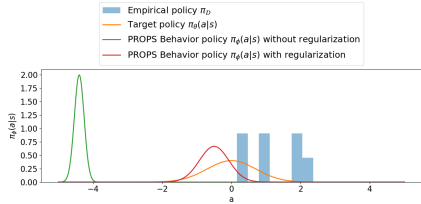


Figure 6: In this example, $\pi(\cdot|s) = \mathcal{N}(0, 1)$. After several visits to s , all sampled actions (blue) satisfy $a > 0$ so that actions $a < 0$ are under-sampled. Without regularization, PROPS will attempt to increase the probabilities of under-sampled action in the tail of target policy distribution (green). The regularization term in the PROPS objective ensures the behavior policy remains close to target policy.

$g(s, \mathbf{a}, \phi, \theta) > 0$	Is the objective clipped?	Return value of min	Gradient
$g(s, \mathbf{a}, \phi, \theta) \in [1 - \epsilon_{\text{PROPS}}, 1 + \epsilon_{\text{PROPS}}]$	No	$-g(s, \mathbf{a}, \phi, \theta)$	$\nabla_{\phi} \hat{\mathcal{L}}_{\text{CLIP}}$
$g(s, \mathbf{a}, \phi, \theta) > 1 + \epsilon_{\text{PROPS}}$	No	$-g(s, \mathbf{a}, \phi, \theta)$	$\nabla_{\phi} \hat{\mathcal{L}}_{\text{CLIP}}$
$g(s, \mathbf{a}, \phi, \theta) < 1 - \epsilon_{\text{PROPS}}$	Yes	$-(1 - \epsilon_{\text{PROPS}})$	$\mathbf{0}$

Table 1: Behavior of PROPS’s clipped surrogate objective (Eq. 5).

C COMPUTING SAMPLING ERROR

We claim that PROPS improves the data efficiency of on-policy learning by reducing sampling error in the agent’s buffer \mathcal{D} with respect to the agent’s current (target) policy. To measure sampling error, we use the KL-divergence $D_{\text{KL}}(\pi_{\mathcal{D}}||\pi_{\theta})$ between the empirical policy $\pi_{\mathcal{D}}$ and the target policy π_{θ} which is the primary metric [Zhong et al. \(2022\)](#) used to show ROS reduces sampling error:

$$D_{\text{KL}}(\pi_{\mathcal{D}}||\pi_{\theta}) = \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\mathcal{D}}(\cdot|s)} \left[\log \left(\frac{\pi_{\mathcal{D}}(a|s)}{\pi_{\theta}(a|s)} \right) \right]. \tag{7}$$

We compute a parametric estimate of $\pi_{\mathcal{D}}$ by maximizing the log-likelihood of \mathcal{D} over the same policy class used for π_{θ} . More concretely, we let θ' be the parameters of neural network with the same architecture as π_{θ} train and then compute:

$$\theta_{\text{MLE}} = \arg \max_{\theta'} \sum_{(s, a) \in \mathcal{D}} \log \pi_{\theta'}(a|s) \tag{8}$$

using stochastic gradient ascent. After computing θ_{MLE} , we then estimate sampling error using the Monte Carlo estimator:

$$D_{\text{KL}}(\pi_{\mathcal{D}}||\pi_{\theta}) \approx \sum_{(s, a) \in \mathcal{D}} (\log \pi_{\theta_{\text{MLE}}}(a|s) - \log \pi_{\theta}(a|s)). \tag{9}$$

D ADDITIONAL EXPERIMENTS

In this appendix, we include additional experiments and ablations.

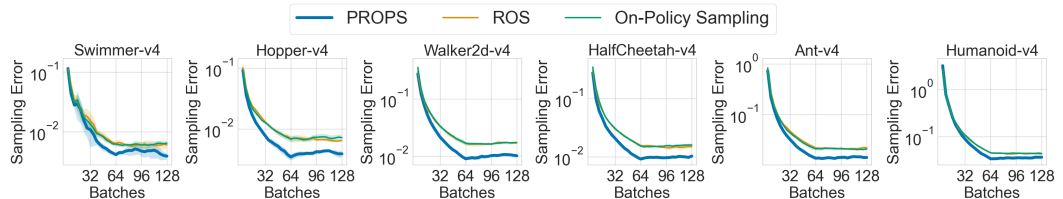


Figure 7: Sampling error with a fixed, expert target policy. Solid curves denote the mean over 5 seeds. Shaded regions denote 95% confidence belts.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

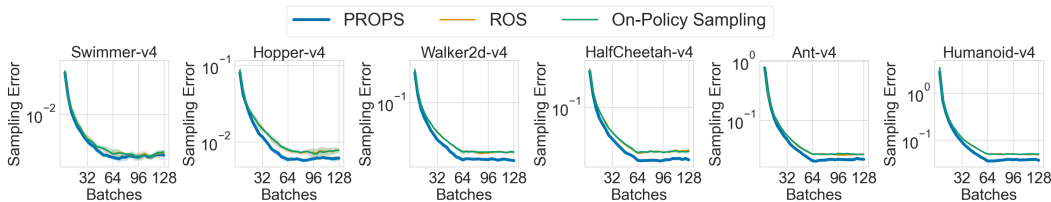


Figure 8: Sampling error with a fixed, randomly initialized target policy. Solid curves denote the mean over 5 seeds. Shaded regions denote 95% confidence belts.

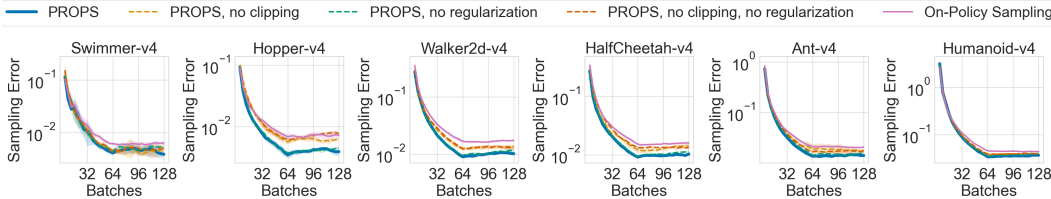


Figure 9: Sampling error ablations with a fixed, expert target policy. Here, “no clipping” refers to setting $\epsilon_{\text{PROPS}} = \infty$, and “no regularization” refers to setting $\lambda = 0$. Solid curves denote the mean over 5 seeds, and shaded regions denote 95% bootstrap confidence intervals.

D.1 CORRECTING SAMPLING ERROR FOR A FIXED TARGET POLICY

In this appendix, we expand upon results presented in Section 6.1 of the main paper and provide additional experiments investigating the degree to which PROPS reduces sampling error with respect to a fixed target policy. We include empirical results for all six MuJoCo benchmark tasks as well as ablation studies investigating the effects of clipping and regularization.

We tune PROPS and ROS using a hyperparameter sweep. For PROPS, we consider learning rates in $\{10^{-3}, 10^{-4}\}$, regularization coefficients $\lambda \in \{0.01, 0.1, 0.3\}$, and PROPS target KLs in $\delta_{\text{PROPS}} \in \{0.05, 0.1\}$. We fix $\epsilon_{\text{PROPS}} = 0.3$ across all experiments. For ROS, we consider learning rates in $\{10^{-3}, 10^{-4}, 10^{-5}\}$. We report results for the hyperparameters yielding the lowest sampling error.

Fig. 7 and 8 show sampling error computed with a fixed expert and randomly initialized target policy, respectively. We see that PROPS achieves lower sampling error than both ROS and on-policy sampling across all tasks. ROS shows little to no improvement over on-policy sampling, highlighting the difficulty of applying ROS to higher dimensional tasks with continuous actions.

Fig. 9 ablates the effects of PROPS’s clipping mechanism and regularization on sampling error reduction. We ablate clipping by setting $\epsilon_{\text{PROPS}} = \infty$, and we ablate regularization by setting $\lambda = 0$. We use a fixed expert target policy and use the same tuning procedure described earlier in this appendix. In all tasks, PROPS achieves higher sampling error without clipping nor regularization than it does with clipping and regularization. However, it nevertheless outperforms on-policy sampling in all tasks except Hopper where it matches the performance of on-policy sampling. Only including regularization slightly decreases sampling error, whereas clipping alone produces sampling error only slightly higher than that achieved by PROPS with both regularization and clipping. These observations indicate that while regularization is helpful, clipping has a stronger effect on sampling error reduction than regularization when the target policy is fixed.

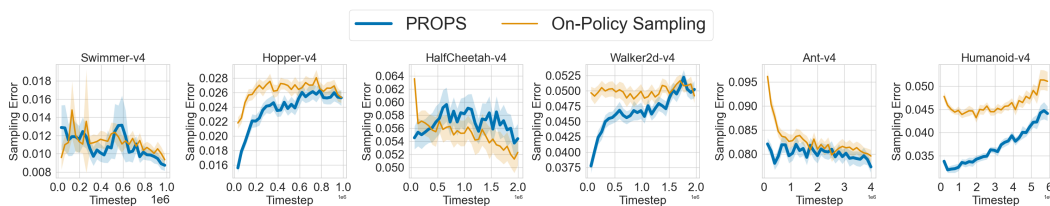


Figure 10: Sampling error throughout RL training. Solid curves denote the mean over 5 seeds. Shaded regions denote 95% confidence belts.

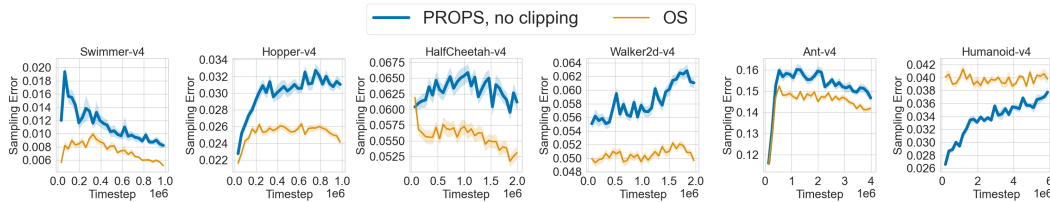


Figure 11: Sampling error throughout RL training without clipping the PROPS objective. Solid curves denote the mean over 5 seeds. Shaded regions denote 95% confidence belts.

D.2 CORRECTING SAMPLING ERROR DURING RL TRAINING

In this appendix, we include additional experiments investigating the degree to which PROPS reduces sampling error during RL training, expanding upon results presented in Section 6.2 of the main paper. We include sampling error curves for all six MuJoCo benchmark tasks and additionally provide ablation studies investigating the effects of clipping and regularization on sampling error reduction and data efficiency in the RL setting.

As shown in Fig 10, PROPS achieves lower sampling error than on-policy sampling throughout training in 5 out of 6 tasks. We observe that PROPS increases sampling error but nevertheless improves data efficiency in HalfCheetah as shown in Fig. 4a. This result likely arises from our tuning procedure in which we selected hyperparameters yielding the largest return. Although lower sampling error intuitively correlates with increased data efficiency, it is nevertheless possible to achieve high return without reducing sampling error.

In our next set of experiments, we ablate the effects of PROPS’s clipping mechanism and regularization on sampling error reduction and data efficiency. We ablate clipping by tuning RL agents with $\epsilon_{\text{PROPS}} = \infty$, and we ablate regularization by tuning RL agents with $\lambda = 0$. Fig. 11 and Fig. 12 show sampling error curves without clipping and without regularization, respectively. Without clipping, PROPS achieves larger sampling than on-policy sampling in all tasks except Humanoid. Without regularization, PROPS achieves larger sampling error in 3 out of 6 tasks. These observations indicate that while clipping and regularization both help reduce sampling during RL training, clipping has a stronger effect on sampling error reduction. As shown in Fig. 13 PROPS data efficiency generally decreases when we remove clipping or regularization.

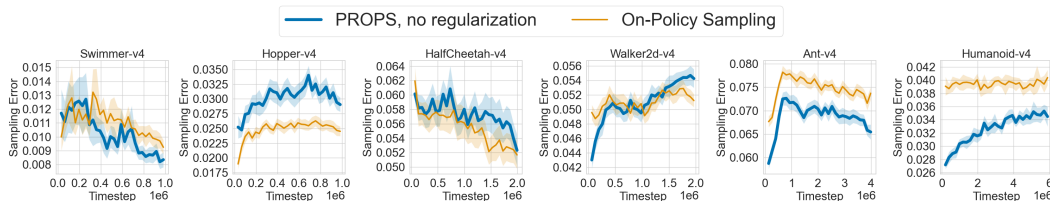


Figure 12: Sampling error throughout RL training without regularizing the PROPS objective. Solid curves denote the mean over 5 seeds. Shaded regions denote 95% confidence belts.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

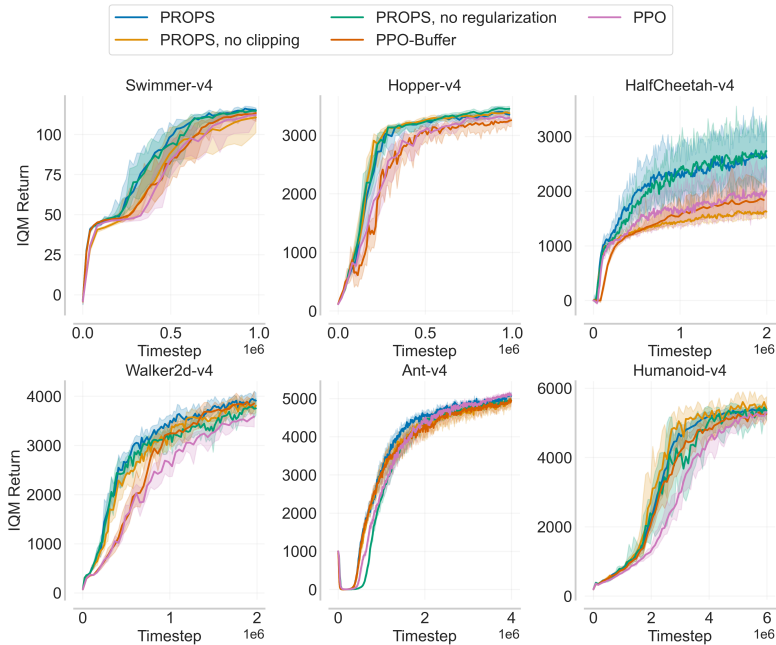


Figure 13: IQM return over 50 seeds of PROPS with and without clipping or regularizing the PROPS objective. Shaded regions denote 95% bootstrap confidence intervals.

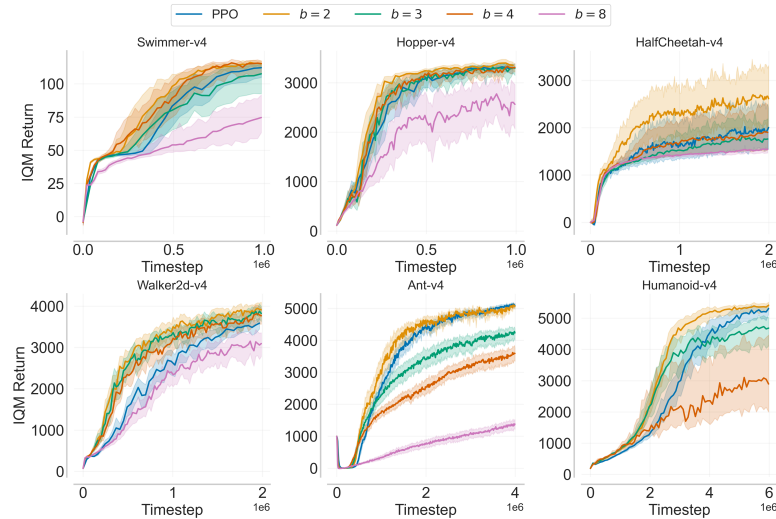


Figure 14: IQM return over 50 seeds for PROPS with different buffer sizes. We exclude $b = 8$ for Humanoid-v4 due to the expense of training and tuning. Shaded regions denote 95% bootstrap confidence intervals.

Lastly, we consider training with larger buffer sizes b in Fig. 14. We find that data efficiency may decrease with a larger buffer size. Intuitively, the more historic data kept around, the more data that must be collected to impact the aggregate data distribution.

D.3 BIAS AND VARIANCE OF PROPS

In Fig. 15, we investigate the bias and variance of the empirical state-action visitation distribution $d^{\mathcal{D}}(s, a)$ under PROPS, ROS, and on-policy sampling. We report the bias and variance averaged over

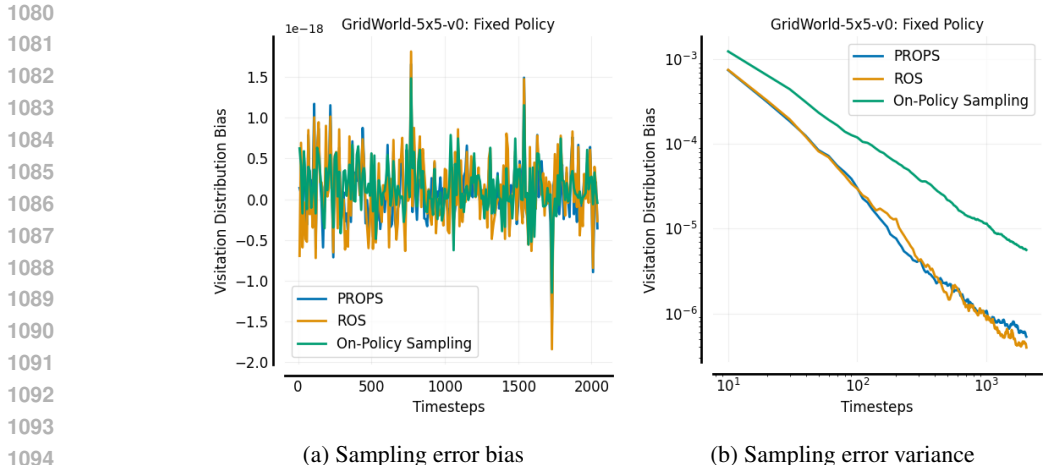


Figure 15: Sampling error bias and variance estimates of different sampling methods. Empirically, PROPS is unbiased and lower variance than on-policy sampling.

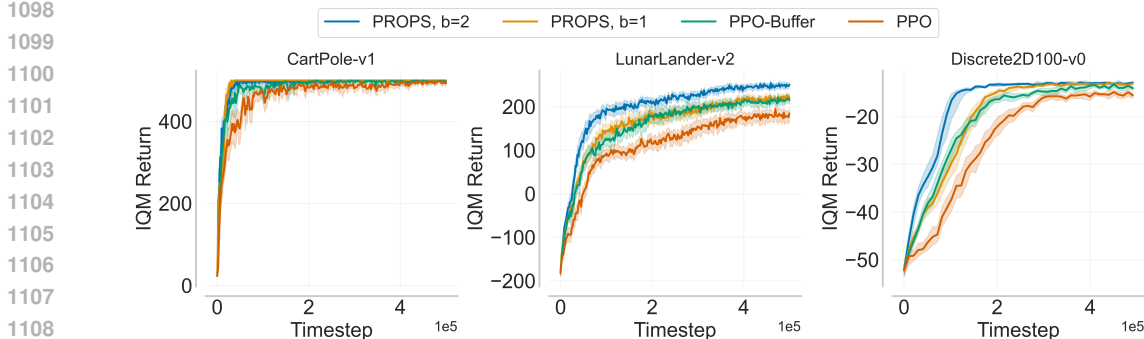


Figure 16: IQM return for discrete action tasks over 50 seeds. Shaded regions denote 95% bootstrap confidence intervals.

all $(s, a) \in \mathcal{S} \times \mathcal{A}$ computed as follows:

$$\text{bias} = \frac{1}{|\mathcal{S} \times \mathcal{A}|} \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} (\mathbb{E} [d^{\mathcal{D}}(s, a)] - d_{\pi_{\theta}}(s, a)) \tag{10}$$

$$\text{variance} = \frac{1}{|\mathcal{S} \times \mathcal{A}|} \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \mathbb{E} \left[(d^{\mathcal{D}}(s, a) - d_{\pi_{\theta}}(s, a))^2 \right] \tag{11}$$

As shown in Fig. 15, the visitation distribution under PROPS and ROS empirical have near zero bias (note that the vertical axis has scale 10^{-18}) and have lower variance than on-policy sampling.

D.4 DISCRETE-ACTION TASKS

We include 3 additional discrete-action domains of varying complexity. The first two are the widely used OpenAI gym domains CartPole-v1 and LunarLander-v2 (Brockman et al., 2016). The third is a 2D navigation task, Discrete2D100-v0, in which the agent must reach a randomly sampled goal. There are 100 actions, each action corresponding to different directions in which the agent can move. From Fig. 16 and 17 we observe that PROPS with $b = 2$ achieves larger returns than PPO and PPO-BUFFER all throughout training in all three tasks. PROPS with $b = 1$ (no historic data) achieves larger returns than PPO all throughout training in all three tasks and even outperforms PPO-BUFFER in CartPole-v1 and Discrete2D100-v0 even though PPO-BUFFER learns from twice as much data. Thus, PROPS can improve data efficiency *without* historic data.

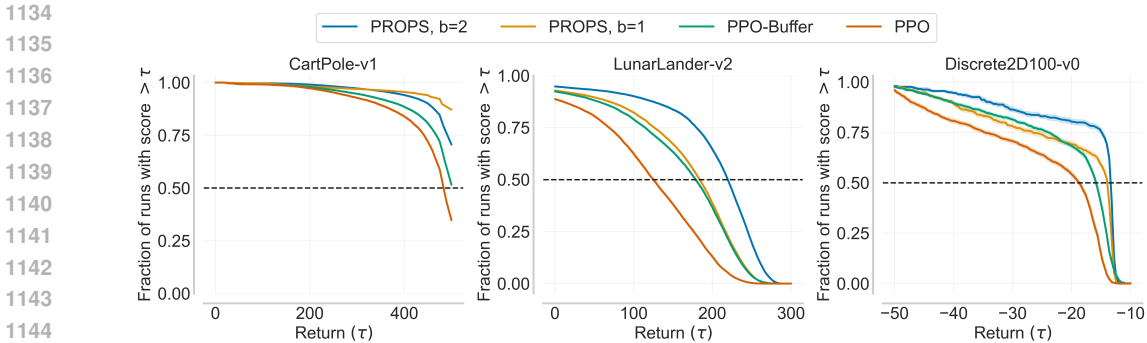


Figure 17: Performance profiles for discrete-action tasks over 50 seeds. Higher values correspond to more reliable convergence to high-return policies. Shaded regions denote 95% bootstrap confidence intervals.

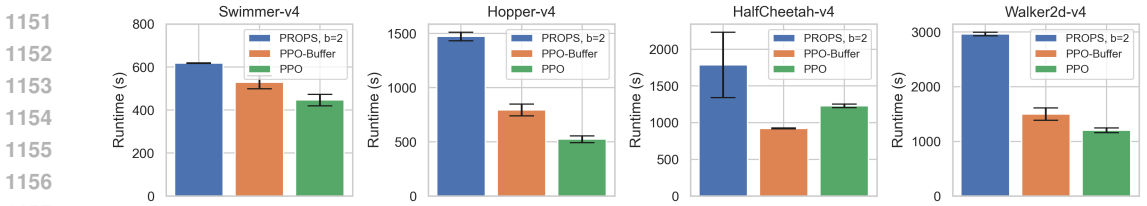


Figure 18: Runtimes for PROPS, PPO-BUFFER, and PPO. We report means and standard errors over 3 independent runs.

D.5 RUNTIME COMPARISONS

Figure 18 shows runtimes for PROPS, PPO-BUFFER, and PPO averaged over 3 runs. We trained all agents on a MacBook Air with an M1 CPU and use the same tuned hyperparameters used throughout the paper. PROPS takes at most twice as long as PPO-BUFFER; intuitively, both PROPS and PPO-BUFFER learn from the same amount of data but PROPS learns two policies.

We note that PPO-BUFFER is faster than PPO is HalfCheetah-v4 because, with our tuned hyperparameters, PPO-BUFFER performs fewer target policy updates than PPO. In particular, PPO-BUFFER is updating its target policy every 4096 steps, whereas PPO is updating the target policy every 1024 steps.

E HYPERPARAMETER TUNING FOR RL TRAINING

For all RL experiments in Section 6.2 and Appendix D.2, we tune PROPS, PPO-BUFFER, and PPO separately using a hyperparameter sweep over parameters listed in Table 2 and fix the hyperparameters in Table 5 across all experiments. Since we consider a wide range of hyperparameter values, we ran 10 independent training runs for each hyperparameter setting. We then performed 50 independent training runs for the hyperparameters settings yielding the largest returns at the end of RL training. We report results for these hyperparameters in the main paper. Fig. 19 shows training curves obtained from a subset of our hyperparameter sweep.

1188	PPO learning rate	$10^{-3}, 10^{-4}$, linearly annealed to 0 over training
1189	PPO batch size n	1024, 2048, 4096, 8192
1190	PROPS learning rate	$10^{-3}, 10^{-4}$ (and 10^{-5} for Swimmer)
1191	PROPS behavior batch size m	256, 512, 1024, 2048, 4096 satisfying $m \leq n$
1192	PROPS KL cutoff δ_{PROPS}	0.03, 0.05, 0.1
1193	PROPS regularizer coefficient λ	0.01, 0.1, 0.3

1194 Table 2: Hyperparameters used in our hyperparameter sweep for RL training.
1195

1196	Environment	Batch Size	Learning Rate
1197	Swimmer-v4	4096	10^{-3}
1198	Hopper-v4	2048	10^{-3}
1199	HalfCheetah-v4	1024	10^{-4}
1200	Walker2d-v4	4096	10^{-4}
1201	Ant-v4	1024	10^{-3}
1202	Humanoid-v4	8192	10^{-4}

1203 Table 3: Tuned PPO hyperparameters
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232

1233	Environment	PPO Batch Size	PPO Learning Rate	PROPS Batch Size	PROPS Learning Rate	PROPS KL Cutoff	PROPS Regularization λ
1234	Swimmer-v4	2048	10^{-3}	1024	10^{-5}	0.03	0.1
1235	Hopper-v4	2048	10^{-3}	256	10^{-3}	0.05	0.3
1236	HalfCheetah-v4	1024	10^{-4}	512	10^{-3}	0.05	0.3
1237	Walker2d-v4	2048	10^{-3}	256	10^{-3}	0.1	0.3
1238	Ant-v4	2048	10^{-4}	256	10^{-3}	0.03	0.1
1239	Humanoid-v4	8192	10^{-4}	256	10^{-4}	0.1	0.1

1240 Table 4: Hyperparameters used in our hyperparameter sweep for RL training.
1241

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

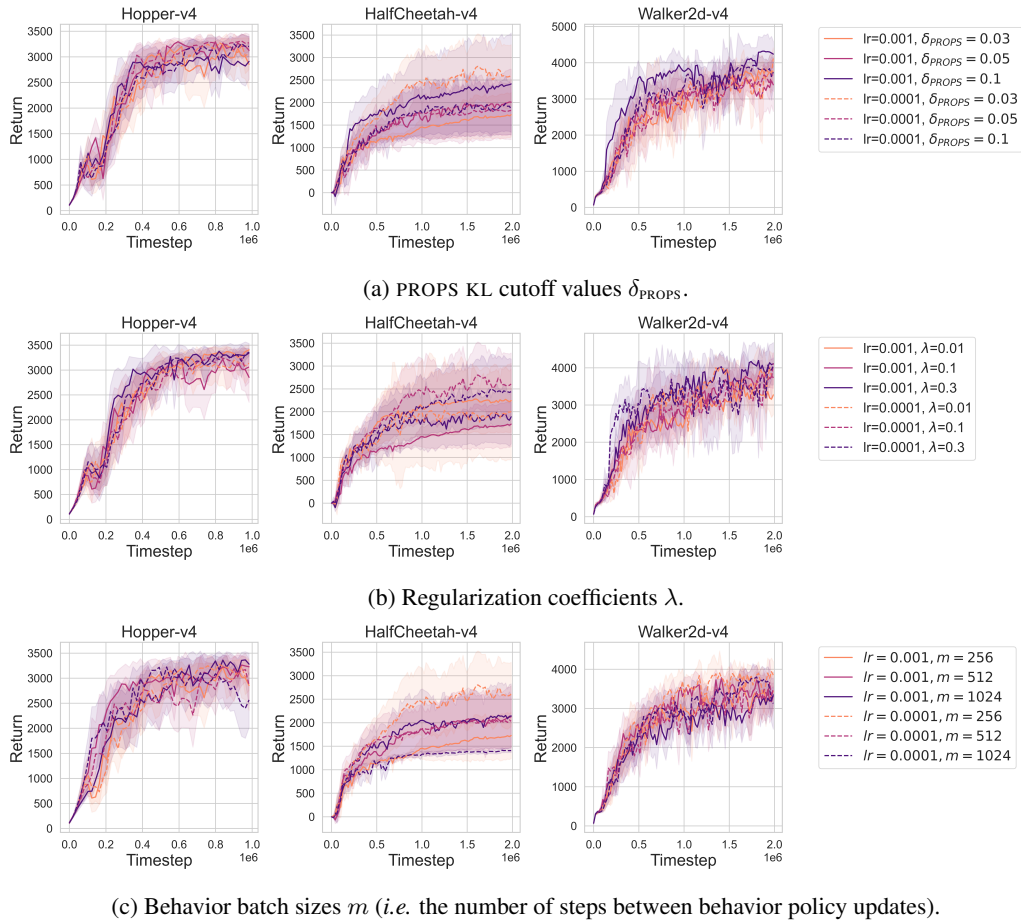


Figure 19: A subset of results obtained from our hyperparameter sweep. Default hyperparameter values are as follows: PROPS KL cutoff $\delta_{PROPS} = 0.03$; regularization coefficient $\lambda = 0.1$; behavior batch size $m = 256$. Darker colors indicate larger hyperparameter values. Solid and dashed lines have the PROPS learning rate set to $1 \cdot 10^{-3}$ and $1 \cdot 10^{-4}$, respectively. Curves denote averages over 10 seeds, and shaded regions denote 95% confidence intervals.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

PPO number of update epochs	10
PROPS number of update epochs	16
Buffer size b	2 target batches (also 3, 4, and 8 in Fig. 14)
PPO minibatch size for PPO update	$bm/16$
PROPS minibatch size for PROPS update	$bm/16$
PPO and PROPS networks	Multi-layer perceptron with hidden layers (64,64)
PPO and PROPS optimizers	Adam (Kingma and Ba, 2015)
PPO discount factor γ	0.99
PPO generalized advantage estimation (GAE)	0.95
PPO advantage normalization	Yes
PPO loss clip coefficient	0.2
PPO entropy coefficient	0.01
PPO value function coefficient	0.5
PPO and PROPS gradient clipping (max gradient norm)	0.5
PPO KL cut-off	0.03
Evaluation frequency	Every 10 target policy updates
Number of evaluation episodes	20

Table 5: Hyperparameters fixed across all experiments. We use the PPO implementation provided by CleanRL (Huang et al., 2022).