

HP3O: HYBRID-POLICY PROXIMAL POLICY OPTIMIZATION WITH BEST TRAJECTORY

Anonymous authors

Paper under double-blind review

ABSTRACT

Proximal policy optimization (PPO) is one of the most popular state-of-the-art on-policy algorithms that has become a standard baseline in modern reinforcement learning with applications in numerous fields. Though it delivers stable performance with theoretical policy improvement guarantees, high variance and high sample complexity still remain critical challenges in on-policy algorithms. To alleviate these issues, we propose Hybrid-Policy Proximal Policy Optimization (HP3O), which utilizes a trajectory replay buffer to make efficient use of trajectories generated by recent policies. Particularly, the buffer applies the "first in, first out" (FIFO) strategy so as to keep only the recent trajectories to attenuate the data distribution drift. A batch consisting of the trajectory with the best return and other randomly sampled ones from the buffer is used for updating the policy networks. The strategy helps the agent to improve its capability on top of the most recent best performance and in turn reduce variance empirically. We theoretically construct the policy improvement guarantees for the proposed algorithm. HP3O is validated and compared against several baseline algorithms using multiple continuous control environments. Our code is available here.

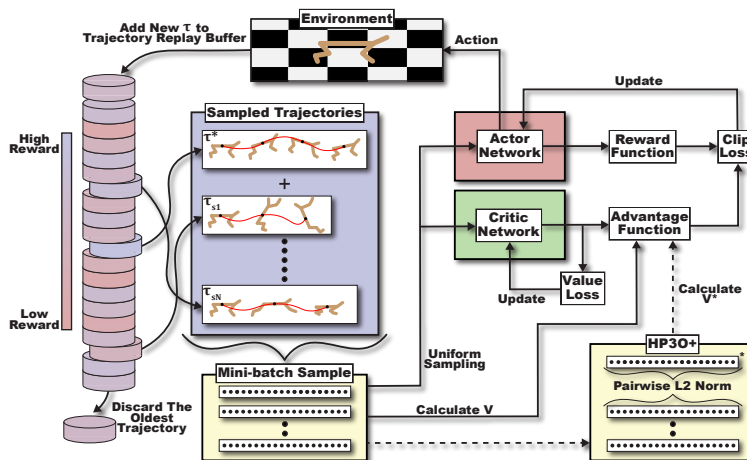


Figure 1: Schematic diagram of HP3O/HP3O+: (left side) the trajectory replay buffer takes a "first in, first out" (FIFO) strategy to keep only recent trajectories; batch consisting of the trajectory with the best return (τ^*) and other randomly sampled ones from the buffer are used for updating the actor/critic networks (*off-policy* approach); (right side) model updating still follows the *on-policy* PPO method, hence, *hybrid-policy* PPO (HP3O); for HP3O+, τ^* is also used to update the advantage function

1 INTRODUCTION

Model-free reinforcement learning Liu et al. (2021) has demonstrated significant success in many different application areas, such as building energy systems Biemann et al. (2021), urban driving Toromanoff et al. (2020); Saxena et al. (2020), radio networks Kaur & Kumar (2020), robotics Polydoros

054 & Nalpantidis (2017), and medical image analysis Hu et al. (2023). In particular, on-policy rein-
 055 forcement learning approaches such as proximal policy optimization (PPO) Schulman et al. (2017);
 056 Chang et al. (2023) provide stable performance along with theoretical policy improvement guarantees
 057 that involve a lower bound Kakade & Langford (2002) on the expected performance loss which
 058 can be approximated using the generated samples from the current policy. These guarantees are
 059 theoretically quite attractive and mathematically elegant, but the requirement of on-policy data and
 060 the high variance nature demands significant data to be collected between every update, inevitably
 061 causing the issue of high sample complexity and the behavior of slow learning.

062 Off-policy algorithms Zanette (2023); Prudencio et al. (2023), on the other hand, alleviate some of
 063 these issues as they can leverage a replay buffer to store samples that enable more efficient policy
 064 updates by reusing these samples. While the off-policy approach leads to better sample efficiency,
 065 it causes another problem called data distribution drift Zhang et al. (2020b); Lesort et al. (2021),
 066 and most studies Lillicrap et al. (2015); Dankwa & Zheng (2019) have just overlooked this issue.
 067 Furthermore, off-policy methods also suffer from high variance and even difficulty in convergence Lyu
 068 et al. (2020) due to the exploration in training. Mitigating this issue Bjorck et al. (2021) still remains
 069 challenging due to the high variations of stored samples in the traditional replay buffer design.
 070 However, it has been receiving considerable attention in recent studies Liu et al. (2020); Xu et al.
 071 (2019). Numerous previous attempts Zhang et al. (2021); Xu et al. (2020); Papini et al. (2018) took
 072 inspiration from supervised learning Wang et al. (2013); Johnson & Zhang (2013) and specifically
 073 made adjustments to the estimation of policy gradients to achieve variance reduction. However, this
 074 involves auxiliary variables and complex estimation techniques, resulting in a more complicated
 075 learning process. Another simple strategy to attenuate high variance is to leverage the advantage
 076 function involving a baseline Jin et al. (2023); Mei et al. (2022); Wu et al. (2018), which can be
 077 estimated by a parameterized model. Nevertheless, when the sampled data from the buffer has a large
 078 distribution drift, learning the parameterized model can be defective, triggering a poor advantage
 value. This naturally leads to the question:

079 *Can we design a hybrid-policy algorithm by assimilating the low sample complexity from off-policy*
 080 *algorithms into on-policy PPO for variance reduction?*

081 **Contributions.** We provide an affirmative answer to the above question. In this work, we blend
 082 off-policy and on-policy approaches to balance the trade-off between sample efficiency and training
 083 stability. Specifically, we focus primarily on mitigating underlying issues of PPO by using a trajectory
 084 replay buffer. In contrast with traditional buffers that keep appending all generated experiences, we
 085 use a "first in, first out" (FIFO) strategy to keep only the recent trajectories to attenuate the data
 086 distribution drift (as shown in Fig. 1). A batch consisting of the trajectory with the best return (a.k.a.,
 087 best trajectory, τ^*) and other randomly sampled ones from the buffer is used for updating the policy
 088 networks. This strategy helps the agent to improve its capability on top of the most recent 'best
 089 performance' and in turn to also reduce variance. Additionally, we define a new baseline which is
 090 estimated from the best trajectory selected from the replay buffer. Such a baseline evaluates how much
 091 better the return is by selecting the present action than the most recent best one, which intuitively
 092 encourages the agent to further improve the performance. More technical detail will be discussed in
 Section 4. Specifically, our contributions are as follows.

- 093 • We propose a novel variant of PPO, called Hybrid-Policy PPO (HP3O), that combines the
 094 advantageous features of on-policy and off-policy techniques to improve sample efficiency
 095 and reduce variance. We also introduce another variant termed HP3O+ that leverages a new
 096 baseline to enhance the model performance. Please see Table 1 for a qualitative comparison
 097 between the proposed and existing methods.
- 098 • We theoretically construct the policy improvement lower bounds for the proposed algorithms.
 099 HP3O provably shows a new lower bound where policies are not temporally correlated,
 100 while HP3O+ induces a value penalty term in the lower bound, which helps reduce the
 101 variance during training.
- 102 • We perform extensive experiments to show the effectiveness of HP3O/HP3O+ across a
 103 few continuous control environments. Empirical evidence demonstrates that our proposed
 104 algorithms are either comparable to or outperform on-policy baselines. Though off-policy
 105 techniques such as soft actor-critic (SAC) may still have better final returns for most tasks,
 106 our hybrid-policy algorithms have significantly more advantages in terms of run time
 107 complexity.

Table 1: Qualitative comparison with PPO and its relevant variants

Method	T.B.	On/off-policy	T.G.
PPO-ClipJin et al. (2023)	X	X	✓
PTR-PPOLiang et al. (2021)	✓	✓	X
GEPPOQueeney et al. (2021)	X	✓	✓
Policy-on-off PPOFakoor et al. (2020)	X	✓	X
P3OChen et al. (2023)	X	X	X
Off-policy PPOMeng et al. (2023)	X	✓	✓
HP3O(+) (ours)	✓	✓	✓

T.B.: trajectory buffer; T.G.: theoretical guarantee.

2 RELATED WORKS

On-policy methods. On-policy algorithms aim at improving the policy performance monotonically between every update. The work Kakade & Langford (2002) developing Conservative Policy Iteration (CPI) for the first time theoretically introduced a policy improvement lower bound that can be approximated by using samples from the present policy. In this regard, trust-region policy optimization (TRPO) Schulman et al. (2015) and PPO have become quite popular baseline algorithms. TRPO solves a trust-region optimization problem to approximately obtain the policy improvement by imposing a Kullback-Leibler (KL) divergence constraint, which requires solving a quadratic programming that may be compute-intensive. On the contrary, PPO achieves a similar objective by adopting a clipping mechanism to constrain the latest policy not to deviate far from the previous one during the update. Their satisfactory performance in different applications Hu et al. (2019); Lele et al. (2020); Zhang et al. (2022); Dutta & Upreti (2022); Bahrpeyma et al. (2023); Nguyen et al. (2024); Zhang et al. (2020a) triggers considerable interest in better understanding these methods Jin et al. (2023) and developing new policy optimization variants Huang et al. (2021). Albeit numerous attempts have been made in the above works, the high sample complexity due to the on-policy behavior of PPO and its variants still obstructs efficient applications to real-world continuous control environments, which demands the connection with off-policy methods.

Off-policy methods. To address the high sample complexity issue in on-policy methods, a common approach is to reuse the samples generated by prior policies, which was devised in Hester et al. (2018); Mnih et al. (2013). Favored off-policy methods such as deep deterministic policy gradient (DDPG) Lillicrap et al. (2015), twin delayed DDPG (TD3) Fujimoto et al. (2018) and soft actor-critic (SAC) Haarnoja et al. (2018) fulfilled this goal by employing a replay buffer to store historical data and sampling from it for computing the policy updates. As mentioned before, such approaches could cause data distribution drift due to the difference between the data distributions of current and prior policies. This work will include an implementation trick to address this issue to a certain extent. Kallus and Uehara developed a statistically efficient off-policy policy gradient (EOPPG) method Kallus & Uehara (2020) and showed that it achieves an asymptotic lower bound that existing off-policy policy gradient approaches failed to attain. Other works such as nonparametric Bellman equation Tosatto et al. (2020) and state distribution correction Kallus & Uehara (2020) were also done with off-policy policy gradient.

Combination of on- and off-policy methods. Making efficient use of on-policy and off-policy schemes is pivotal to designing better model-free reinforcement learning approaches. An early work merged them together to come up with the interpolated policy gradient Gu et al. (2017) for improving sample efficiency. Another work Fakoor et al. (2020) developed Policy-on-off PPO to interleave off-policy updates with on-policy updates, which controlled the distance between the behavior and target policies without introducing any additional hyperparameters. Specifically, they utilized a complex gradient estimate to account for on-policy and off-policy behaviors, which may result in larger computational complexity in low-sample scenarios. To compensate data inefficiency, Liang et al. Liang et al. (2021) incorporated prioritized experience replay into PPO by proposing a truncated importance weight method to overcome the high variance and designing a policy improvement loss function for PPO under off-policy conditions. A more recent work Chen et al. (2023) probed the insufficiency of PPO under an off-policy measure and explored in a much larger policy space to maximize the CPI objective. The most related work to ours is Queeney et al. (2021), where the

162 authors proposed a generalized PPO with off-policy data from prior policies and derived a generalized
 163 policy improvement lower bound. They utilized directly the past trajectories right before the present
 164 one instead of a replay buffer, which still maintains a weakly on-policy behavior. However, their
 165 method may suffer from poor performance in sparse reward environments.

167 3 PROBLEM FORMULATION AND PRELIMINARY

169 **Markov decision process.** In this context, we consider an infinite-horizon Markov Decision Process
 170 (MDP) with discounted reward defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r, \rho_0, \gamma)$, where \mathcal{S} indicates the
 171 set of states, \mathcal{A} signifies the set of actions, $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is the transition probability function,
 172 $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, ρ_0 is the initial state distribution of environment, and
 173 γ is the discount factor. In this study, the agent’s policy is a stochastic mapping represented by
 174 $\pi : \mathcal{S} \rightarrow \mathcal{A}$. Reinforcement learning aims at choosing a policy that is able to maximize the
 175 expected discounted cumulative rewards $J(\pi) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$, where $\tau \sim \pi$ indicates
 176 a trajectory sampled according to $s_0 \sim \rho_0$, $a_t \sim \pi(\cdot|s_t)$, and $s_{t+1} \sim p(\cdot|s_t, a_t)$. We denote by
 177 $d^\pi(s)$ a normalized discounted state visitation distribution such that $d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t =$
 178 $s | \rho_0, \pi, p)$. Hence, the corresponding normalized discounted state-action visitation distribution can be
 179 expressed as $d^\pi(s, a) = d^\pi(s) \pi(s, a)$. Additionally, we define the state value function of the policy
 180 π as $V^\pi(s) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s]$, the state-action value function, i.e., Q -function,
 181 as $Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a]$, and the critical advantage function as
 182 $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$.

183 **Policy improvement guarantee.** The foundation of numerous on-policy policy optimization al-
 184 gorithms is built upon a classic policy improvement lower bound originally established in Kakade
 185 & Langford (2002). With different scenarios Schulman et al. (2015); Achiam et al. (2017); Dai &
 186 Gluzman (2021), the lower bound was refined to reflect diverse policy improvements, which can be
 187 estimated by using the samples generated from the latest policy. For completeness, we present in
 188 Lemma 1 the policy improvement lower bound from Achiam et al. (2017).

189 **Lemma 1.** (Corollary 1 in Achiam et al. (2017)) Suppose that the current time step is k and that the
 190 corresponding policy is π_k . For any future policy π , the following relationship holds true:

$$191 J(\pi) - J(\pi_k) \geq \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d^{\pi_k}} \left[\frac{\pi(a|s)}{\pi_k(a|s)} A^{\pi_k}(s, a) \right] - \frac{2\gamma C_{\pi_k}^\pi}{(1 - \gamma)^2} \mathbb{E}_{(s,a) \sim d^{\pi_k}} [\delta(\pi, \pi_k)(s)], \quad (1)$$

192 where $C_{\pi_k}^\pi = \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]|$ and $\delta(\pi, \pi_k)(s)$ is the total variation distance between
 193 the distributions $\pi(\cdot|s)$ and $\pi_k(\cdot|s)$.

194 Lemma 1 implies that the policy improvement lower bound consists of the surrogate objective loss
 195 and the penalty term, which can be maximized by choosing a certain new policy π_{k+1} to guarantee
 196 the policy improvement. However, directly maximizing such a lower bound could be computationally
 197 intractable if the next policy π_{k+1} deviates far from the current one. Unless additional constraint
 198 is imposed such as a trust region in TRPO Schulman et al. (2015), which unfortunately requires a
 199 complex second-order method to solve the optimization problem. Hence, PPO developed a simple
 200 yet effective heuristic for achieving this.

202 **Proximal policy optimization.** PPO has become a default baseline in a variety of applications, as
 203 mentioned above. It is favored because of its strong performance and simple implementation with
 204 sound theoretical motivation given by the policy improvement lower bound. Intuitively, PPO attempts
 205 to constrain the new policy close to the present one with a *clipping* heuristic, which results in the
 206 most popular variant, PPO-clip Jin et al. (2023). Particularly, the following objective is solved at
 207 every policy update:

$$208 \mathcal{L}_k^{clip}(\pi) = \mathbb{E}_{(s,a) \sim d^{\pi_k}} \left[\min \left(\frac{\pi(a|s)}{\pi_k(a|s)} A^{\pi_k}(s, a), \text{clip} \left(\frac{\pi(a|s)}{\pi_k(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A^{\pi_k}(s, a) \right) \right], \quad (2)$$

209 where $\text{clip}(a, b, c) = \min(\max(a, b), c)$. The clipping function plays a critical role in this objective
 210 as it consistently enforces the probability ratio between the current and next policies in a reasonable
 211 range between $[1 - \epsilon, 1 + \epsilon]$. The outer minimization in Eq. 2 provides the lower bound guarantee for
 212 the surrogate loss in Eq. 1. In practice, one can set a small learning rate and a large number of time
 213 steps to generate sufficient samples to allow PPO to perform stably and approximate Eq. 2. However,
 214 due to its on-policy approach, high variance is a significant issue such that an extremely large number
 215 of samples may be required in some scenarios to make sure the empirical objective is able to precisely
 estimate the true objective in Eq. 2, which naturally causes the high sample complexity issue. This

motivates us to leverage off-policy techniques to alleviate such an issue, while keeping the theoretical policy improvement.

4 HYBRID-POLICY PPO (HP3O)

To achieve better sample efficiency of PPO, historical samples generated by previous policies are reused for policy updates, as done in off-policy algorithms. This inevitably results in a distribution drift between policies, which essentially disproves the policy improvement lower bound in Lemma 1. In this context, to fix this issue, we will extend Lemma 1 to assimilate off-policy samples in a principled manner to derive a new policy improvement lower bound that works for our proposed algorithm, HP3O. HP3O (and its variant HP3O+) takes a *hybrid* approach that effectively synthesizes on-policy trajectory-wise policy updates and off-policy trajectory replay buffers. Algorithm 1 shows

Algorithm 1 HP3O(+)

```

1: Input: initializations of  $\theta_0$ ,  $\phi_0$ , and trajectory replay buffer  $R$ , the number of episodes  $K$ , the
2: number of time steps in each episode  $T$ , the number of epochs for updates  $E$ 
3: for  $k = 1, 2, \dots, K$  do
4:   Run policy  $\pi_{\theta_k}$  to generate a trajectory  $\tau = (s_0, a_0, r_1, s_1, \dots, s_{T-1}, a_{T-1}, r_T)$ 
5:   Append  $\tau$  to  $R$  and discard the oldest one  $\tau^-$  ▷ FIFO strategy
6:   Sample a random minibatch  $\mathcal{B}$  from the trajectory replay buffer  $R$ 
7:   Select the best action trajectory  $\tau_k^*$  from the trajectory replay buffer and add it to  $\mathcal{B}$ 
8:   for each trajectory  $j = 1, 2, \dots, |\mathcal{B}|$ : do
9:     for  $t = 0, 1, \dots, T - 1$  do
10:       $G_t^j = \sum_{l=t+1}^T \gamma^{l-t-1} r_l^j$ 
11:    end for
12:  end for
13:  Compute advantage estimates  $\hat{A}_t^{\pi_k} = G_t - V_\phi(s_t)$  ▷ HP3O
14:  Compute  $V^{\tau_k^*}(s_t)$  using  $\tau_k^*$  and advantage estimates  $\hat{A}_t^{\pi_k} = G_t - V^{\tau_k^*}(s_t)$  ▷ HP3O+
15:  for each epoch  $e = 1, 2, \dots, E$  do
16:    Compute the clipping loss Eq. 2
17:    Compute the mean square loss  $\mathcal{L}^V(\phi) = -\frac{1}{T} \sum_{t=0}^{T-1} (G_t - V_\phi(s_t))^2$ 
18:    Update  $\pi_{\theta_k}$  with  $\nabla_\theta \mathcal{L}^{clip}(\theta)$  by Adam
19:    Update  $V_{\phi_k}$  with  $\nabla_\phi \mathcal{L}^V(\phi)$  by Adam
20:  end for
21: return  $\pi_{\theta_K}$  and  $V_{\phi_K}$ 

```

the algorithm framework for HP3O and HP3O+ (blue line represents the only difference for HP3O+). We denote the actor and critic by $\theta \in \mathbb{R}^m$ and $\phi \in \mathbb{R}^n$ respectively such that the parameterized policy function is π_θ and the parameterized value function is $V_\phi = \mathbb{E}_{\tau \sim \pi_\theta} [\sum_{l=t}^T \gamma^{l-t} r(s_l, a_l) | s_t]$. Denote by $\tau_k^* = \operatorname{argmax}_{\tau \in R} \sum_{t=0}^T \gamma^t r(s_t, a_t)$ the best action trajectory selected from the replay buffer R at the current episode k .

In most existing off-policy algorithms, the size of the replay buffer is fixed with a large number to ensure that a diverse set of experiences is captured. With this approach, though the random minibatch sampling allows the agent to learn from past experience, a large-size replay buffer may cause significant data distribution drifts. Additionally, a large replay buffer means that it takes more time for the buffer to fill up, especially in environments requiring extensive exploration. Hence, we apply the FIFO strategy and discard old trajectories empirically to attenuate the issue (Line 4 in Algorithm 1). The recently proposed off-policy PPO Meng et al. (2023) indeed uses off-policy data, but it does not employ a trajectory buffer as we do. In our approach, the trajectory buffer is an essential component because it allows us to store and process complete sequences of state-action pairs (trajectories) rather than isolated transitions. This will preserve the temporal coherence and enhance stability. Line 5 is to sample from the trajectory replay buffer R , which is different from the reuse of N samples generated from prior policies in Queeney et al. (2021), where the past immediate sample trajectories were used without random sampling. We note that a replay buffer in the proposed algorithm enhances the agent’s performance by providing access to a more diverse

set of experiences and highlighting the most impactful trajectories. Line 6 signifies the core part of HP3O as the best action trajectory τ_k^* indicates the best return starting from state s_t within the buffer. Line 7 through Line 12 calculate the rewards to go for each time step t in each trajectory and obtain the total reward to go at each time step over all trajectories. One may wonder how to calculate the return G_t if trajectories have varying lengths in some environments. In this work, we store different lengths of trajectories directly in the buffer and do not pad them. This approach preserves the natural variation in trajectory lengths that can occur in different environments. Although the length differ, we still compare the returns of these trajectories to identify the best one while ensuring the comparison remains consistent and fair. Particularly, line 13 is a key step in the proposed HP3O+. $V^{\tau_k^*}(s_t) = \mathbb{E}_{\tau_k^* \sim \pi_k} [\sum_{l=t}^T \gamma^{l-t} r(s_l, a_l) | s_t]$ induced by the current best action trajectory τ_k^* sets the best state value among all trajectories from R . $\hat{A}_t^{\pi_k}$ in Line 13 signifies how much better the return G_t is by taking action a_t than the best value we have obtained most recently. Intuitively, this "encourages" the agent to improve its performance in the next step on top of $V^{\tau_k^*}(s_t)$. While $V^{\tau_k^*}(s_t)$ can be theoretically calculated as above, in practice, to make sure that there always exists a best value for use, $V^{\tau_k^*}(s_t)$ is calculated by using a norm distance between the current trajectory and best trajectories to ensure $V^{\tau_k^*}$ has the best return since s_t . If the reward to go from s_t in the best trajectory is lesser, the current trajectory is used to replace the best one for $V^{\tau_k^*}(s_t)$ calculation. Please see the Appendix for more details about the data structures of the proposed algorithms.

Remark 1. *We remark on the sampling method adopted in this work to obtain the trajectories apart from the best trajectory for update. We begin by randomly sampling a set of trajectories from our trajectory buffer. This set is specifically designed to include the best action trajectory, with the remaining trajectories selected randomly from the buffer. From the set of trajectories obtained by random sampling, we then apply uniform sampling. The resulting minibatch is used for training. This approach balances leveraging high-performing trajectories while maintaining exploration across the broader trajectory space, helping to reduce the risk of overfitting. However, we recognize that assigning a score to trajectories based on the loss function could offer additional benefits. Prioritizing trajectories Hou et al. (2017) that result in higher losses could help the agent focus on challenging experiences, potentially improving learning efficiency by addressing areas where the policy requires more refinement. This could also help in stabilizing training by emphasizing learning from mistakes, thereby potentially reducing the variance in policy updates. In fact, integrating a prioritized experience replay (PER) strategy could be a promising direction for future work.*

5 THEORETICAL ANALYSIS

This section presents a theoretical analysis of the proposed HP3O and HP3O+. We first derive a new policy improvement lower bound for HP3O and then present a different bound for HP3O+ to indicate the value penalty term. All proofs are deferred to the Appendix. To incorporate prior policies in the policy improvement lower bound, we need to extend the conclusion in Lemma 1, which quantifies the improvement for two consecutive policies. In Queeney et al. (2021), policies prior to the present policy π_k in chronological order were used. However, in our study, this order has been broken due to the random sampling from the replay buffer, which motivates us to derive a relationship among the current, future, and prior policies independent of the chronological order. Before the main result, we first present an auxiliary technical lemma.

Lemma 2. *Consider a current policy π_k , and any reference policy π_r . For any future policy π ,*

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi_r}} \left[\frac{\pi(a|s)}{\pi_r(a|s)} A^{\pi_k}(s, a) \right] - \frac{2\gamma C_{\pi_k}^{\pi}}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)], \quad (3)$$

where $C_{\pi_k}^{\pi}$ and $\delta(\pi, \pi_r)(s)$ are defined as in Lemma 1.

Remark 2. *Lemma 2 implies that now the visitation distribution, the probability ratio of the surrogate objective, and the maximum value of the total variation distance depend on the reference policy π_r , which essentially extends Lemma 1 to a more generalized case. However, the improvement is still for the two consecutive policies π_k and π as the advantage function in the surrogate objective and $C_{\pi_k}^{\pi}$ rely on the latest policy π_k . Lemma 2 does not necessarily require π_r to be the last policy prior to π_k as in Queeney et al. (2021), which paves the way for establishing the policy improvement for $|\mathcal{B}|$ prior policies sampled randomly from the replay buffer R .*

Theorem 1. *Consider prior policies $|\mathcal{B}|$ randomly sampled from the replay buffer R with indices $i = 0, 1, \dots, |\mathcal{B}| - 1$. For any distribution $v = [v_1, v_2, \dots, v_{|\mathcal{B}|}]$ over the $|\mathcal{B}|$ prior policies, and any*

324 future policy π generated by HP3O in Algorithm 1, the following relationship holds true

$$325 \quad J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{i \sim v} [\mathbb{E}_{(s,a) \sim d^{\pi_i}} [\frac{\pi(a|s)}{\pi_i(a|s)} A^{\pi_k}(s,a)]] - \frac{\gamma C_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2}, \quad (4)$$

327 where $C_{\pi_k}^{\pi}$ is defined as in Lemma 1.

328 **Remark 3.** It is observed that the conclusion from Theorem 1 is similar to one of the main results
329 in Queeney et al. (2021). The significant difference is that π_i is not the same as π_{k-i} in Queeney et al.
330 (2021). It is technically attributed to Lemma 2, where the reference policy π_r may not have a close
331 temporal relationship with π_k . Also, the advantage function has not been changed yet. Empirically
332 speaking, for each minibatch \mathcal{B} , we have added the best trajectory in it, which essentially expedites
333 the learning process. Additionally, Theorem 1 has an extra expectation operator over multiple
334 trajectories on the first term of the right side in Eq. 4, leading to the smaller variance, compared to
335 only one trajectory in Lemma 1. We would also like to point out that Theorem 1 shows the policy
336 improvement lower bound by sampling a mini-batch of trajectories associated with prior policies
337 from the buffer, which is consistent with what has been done in Algorithm 1. In HP3O+, we use it as
338 a baseline to replace $V_{\phi}(s)$ and have surprisingly found that this leads to an extra term that penalizes
339 the state value to reduce the variance.

340 We first define $\hat{A}^{\pi}(s,a) = Q^{\pi}(s,a) - V^{\pi^*}(s)$ and $G^{\pi}(s) = V^{\pi^*}(s) - V^{\pi}(s)$. It is immediately
341 obtained that $A^{\pi}(s,a) = \hat{A}^{\pi}(s,a) + G^{\pi}(s)$. Hence, if we utilize the state value induced by the
342 best trajectory at the moment as the baseline, there exists a *value gap* $G^{\pi}(s)$ between $A^{\pi}(s,a)$ and
343 $\hat{A}^{\pi}(s,a)$. One may argue that the advantage $\hat{A}^{\pi}(s,a)$ is negative all the time, which implies the
344 present action is not favorable such that the new policy should be changed to yield a lower probability
345 for the current action and state. However, this is not always true as $V^{\pi^*}(s)$ is not the *globally*
346 optimal value, while it is approximately the optimal value up to the current time step over the last
347 $|\mathcal{B}|$ episodes. The motivation behind $\hat{A}^{\pi}(s,a)$ is that the new baseline $V^{\pi^*}(s)$ becomes the driving
348 force to facilitate the performance improvement between every update. We are now ready to state the
349 policy improvement lower bound with the new baseline as follows.

350 **Lemma 3.** Consider a current policy π_k , and any reference policy π_r . For any future policy π ,

$$351 \quad J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi_r}} [\frac{\pi(a|s)}{\pi_r(a|s)} \hat{A}^{\pi_k}(s,a)] - \frac{2\gamma \hat{C}_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)] \\ 352 \quad - \frac{2\gamma C^{\pi_k}}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)], \quad (5)$$

353 where $\hat{C}_{\pi_k}^{\pi} = \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} [\hat{A}^{\pi_k}(s,a)]|$, $\delta(\pi, \pi_r)(s)$ is defined as in Lemma 1, $C^{\pi_k} =$
354 $\max_{s \in \mathcal{S}} |V^{\pi_k^*}(s) - V^{\pi_k}(s)|$.

355 With Lemma 3 in hand, we have another main result in the following.

356 **Theorem 2.** Consider prior policies $|\mathcal{B}|$ randomly sampled from the replay buffer R with indices
357 $i = 0, 1, \dots, |\mathcal{B}| - 1$. For any distribution $v = [v_1, v_2, \dots, v_{|\mathcal{B}|}]$ over the $|\mathcal{B}|$ prior policies, and any
358 future policy π generated by HP3O+ in Algorithm 1, the following relationship holds true

$$359 \quad J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{i \sim v} [\mathbb{E}_{(s,a) \sim d^{\pi_i}} [\frac{\pi(a|s)}{\pi_i(a|s)} \hat{A}^{\pi_k}(s,a)]] - \frac{\gamma \hat{C}_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2} - \frac{\gamma C^{\pi_k} \epsilon}{(1-\gamma)^2}, \quad (6)$$

360 where $\hat{C}_{\pi_k}^{\pi}$ and C^{π_k} are defined as in Lemma 3.

361 **Remark 4.** Theorem 2 describes the policy improvement lower bound for HP3O+, which provides
362 the theoretical guarantees when reusing trajectories generated by prior policies rigorously. The
363 extra term on the right-hand side $\frac{\gamma \hat{C}_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2}$ in the above inequality is not the penalty term between two
364 policies, while it is a *value gap* between the current state value and the most recent best value. As
365 $V^{\pi_k^*}(s)$ is time-varying, this acts as a "guide" to the current one V^{π_k} not deviating too far away from
366 $V^{\pi_k^*}(s)$. Equivalently, the term $\frac{\gamma C^{\pi_k} \epsilon}{(1-\gamma)^2}$ can be regarded as a regularization from the critic network,
367 which assists in enhancing the overall agent performance and reducing the variance. We also include
368 some technical discussion regarding whether our approach will cause overfitting and the adoption of
369 the worst trajectories in Appendix A.2 and A.3.

370 **Remark 5.** The proposed HP3O algorithm and its variant have resorted to data randomly sampled
371 from multiple policies in the training batch \mathcal{B} that is prior to π_k for the policy update. Thus, there
372 exist multiple updates compared to the vanilla PPO, which only makes one policy update from π_k
373

to π_{k+1} . In this study, we aim to show how the off-policy sample reuse significantly affects the original sample efficiency PPO has. Though the direct sample complexity improvement analysis can be significantly beneficial to provide a solid theoretical foundation for the proposed algorithms, a thorough investigation of this aspect is out of the scope of this study. For instance, to arrive at an ε -optimality for policy gradient-based algorithms, a few works Zhong & Zhang (2024); Zanette et al. (2021); Dai et al. (2023); Sherman et al. (2023) have revealed the exact complexity with respect to ε , but only for MDPs with linear function approximation. The exact sample complexity analysis for the off-policy PPO algorithm with nonlinear function approximation still remains extremely challenging and requires a substantial amount of non-trivial efforts. Therefore, in this paper, we instead disclose the impact of off-policy sample reuse on the tradeoff between sample efficiency and learning stability. Please see Appendix A.4 for more details. Additionally, we also present a theoretical result in Appendix A.5 to reveal that HP3O+ increases updates in the total variational distance of the policy throughout training, given the same sample size, when it is compared to HP3O.

6 EXPERIMENTS

The experimental evaluation aims to understand how the sample complexity and stability of our proposed algorithms compare with existing baseline on-policy and off-policy learning algorithms. Concretely, we conduct the comparison between our methods and prior approaches across challenging continuous control environments from the Gymnasium benchmark suite Brockman et al. (2016). While easy control tasks can be solved by various algorithms, the more complex tasks are typically sample intensive with on-policy algorithms Schulman et al. (2017). Additionally, the high variance of the algorithms negatively impacts stability and convergence. Furthermore, though some off-policy algorithms enjoy high sample efficiency, the actual run time can be impractically large, which impedes its applications to real-world tasks. As our proposed hybrid-policy learning algorithms are developed on top of PPO, we mainly compare our methods to PPO, another popular on-policy method A2C Peng et al. (2018), and three other relevant off-policy PPO approaches, including P3O Chen et al. (2023) (a modification of PPO to leverage both on- and off-policy principles), GEPPPO Queeney et al. (2021), and Off-policy PPO (abbreviated as OffPolicy) Meng et al. (2023). We acknowledge that SAC, a fully off-policy algorithm, may achieve comparatively higher returns in most of the continuous control problems at the expense of much longer training time and with careful hyperparameter tuning. Hence, we also compare with SAC in terms of variance reduction and run time complexity. As shown in Table 1, there are other off-policy versions of PPO, such as Policy-on-off PPO Fakoor et al. (2020). However, the corresponding code base lacks a complete implementation to reproduce their results, which is evident in their code where the actor head for Mujoco environments is not implemented. Moreover, making the code functional for our purpose would require extensive effort, as it is built on MXNet, a deprecated open-source project. The above limitations have prevented us from performing head-to-head comparisons. More details about hyperparameter settings are deferred to the Appendix.

6.1 COMPARATIVE EVALUATION

Figure 2 shows the total average return during training for A2C, PPO, P3O, GEPPPO, OffPolicy, HP3O, and HP3O+. Each experiment includes five different runs with various random seeds. The solid curves indicate the mean, while the shaded areas represent the standard deviation over the five runs. Clearly, the results show that, overall, both HP3O+ and HP3O are comparable to or outperform all baselines across diverse tasks with smaller variances, which supports our theoretical claims. For instance, in the HalfCheetah environment, our methods demonstrate a sharper average slope compared to the baseline, particularly in the later stages of training, where other baselines show a more flattened curve. This indicates that our method continues to learn effectively with fewer samples. In the Hopper environment, P3O performs slightly better than HP3O but at the cost of extremely large reward variance, indicating an unstable training process. However, HP3O+ significantly dominates in the latter phase with a much smaller variance. In the Swimmer environment, while A2C and P3O learn slowly and make almost no progress, HP3O and HP3O+ achieve the similarly highest reward with very low variance, as suggested by Remark 3. Notably, OffPolicy ranks second in terms of performance, but with the cost of extremely high variance. Additionally, OffPolicy shows notably strong performance in the Walker environment. This is primarily attributed to the adoption of a new clipped surrogate that iteratively resorts to off-policy data to progress during training. Generally, our proposed methods learn more stably than all baselines by dequeuing the buffer to suppress the instability caused by data distribution drift in most environments. Overall, HP3O+ excels HP3O in most environments, with also variance reduction particularly in the latter training phase. As the learning trajectories are always around the best trajectory from the buffer. Essentially, the empirical

evidence supports our theoretical results in Theorem 2 and Theorem 5, which show that HP3O+ enables larger updates in the total variational distance of the policy, given the same number of changes to the policy. Additional results are included in the Appendix, including Table 2 to showcase rewards at or close to the converged stage.

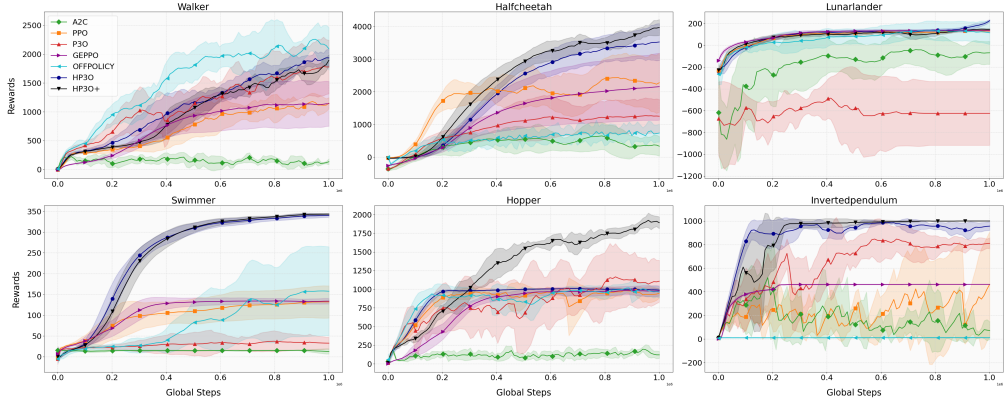
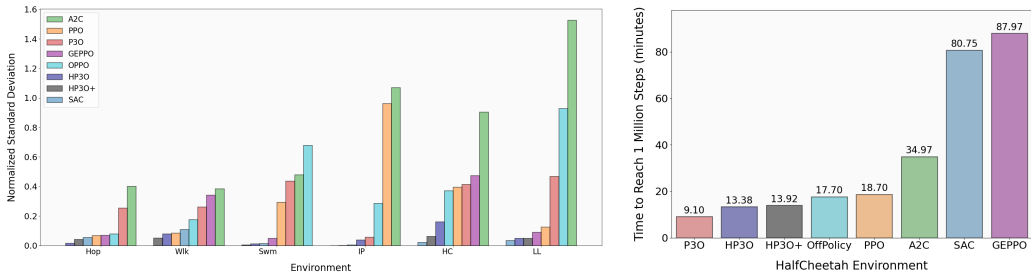


Figure 2: Training curves (over 1M steps) on continuous control benchmarks. HP3O+ (black) performs consistently across all tasks and is comparable to or outperforming other baseline methods.



(a) Normalized Standard Deviation among different methods for various environments. (b) Runtime for HalfCheetah Environment among different methods

Figure 3: Comparison of Normalized Standard Deviation and Runtime for 1 million steps.

6.2 ABLATION STUDY

The experimental results in the previous section imply that algorithms based on the hybrid-policy approach can outperform the conventional on-policy methods on challenging control tasks. In this section, we further compare all policy optimization algorithms to SAC for variance reduction and run time complexity. We also inspect the robustness of the algorithms against variations of trajectories.

Variance. Figure 3a shows the comparison of the relative standard deviation of the ultimate average return (at 1M steps) for different algorithms. It suggests that, on average, HP3O+ achieves the lowest relative standard deviation (which is the ratio of the standard deviation to the average reward over five runs at the last step). This implies that hybrid-policy algorithms have more advantages in regularizing the learning process to maintain stability compared to typical on-policy algorithms. Intuitively, as the policy and environment change over time, the use of replay buffers helps mitigate this issue by providing a more stationary training dataset. The buffer contains a mix of experiences collected under different policies, instead of the only current policy from PPO, which helps in reducing the variance in updates. SAC attains a relatively small standard deviation according to Figure 3a (also, on average, the maximum reward reported in the Appendix). This is not surprising since the maximum entropy principle can significantly help meaningful exploration to achieve the highest return. However, this comes at the cost of runtime complexity.

Run time complexity. As shown in Figure 3b, the run time for all algorithms is presented (all methods are implemented with the same hardware). Both GEPPPO and SAC require much more run time to explore and then converge, which may impede its applications to solving real-world problems. P3O achieves the lowest run time complexity while performing worse than HP3O and

486 **HP3O+**. However, our proposed approaches take approximately the same training time as PPO but
 487 with higher sample efficiency, as shown in Figure 2. Thus, HP3O/HP3O+ are able to achieve a
 488 desirable trade-off in practice between sample efficiency and computational time. These experiments
 489 used a local machine with an NVIDIA RTX 4090. [Additional results regarding wall-clock time for](#)
 490 [diverse methods to reach a certain reward are included in Appendix A.13.](#)

491 **Robustness.** We also compute the *explained variance* LaHuis et al. (2014) for all algorithms under
 492 consideration for evaluating robustness. Please check the Appendix A.7 for more details about this
 493 metric. Intuitively, it quantifies how good a model is to explain the variations in the data. Therefore,
 494 the higher the explained variance of a model, the more the model is able to explain the variations in
 495 trajectories. Essentially, the data in this work are trajectories produced by different policies, leading to
 496 a data distribution drift. Therefore, explained variance can, to some extent, be viewed as an indicator
 497 of how well an algorithm is robust against the data distribution drift. Figure 4 shows the explained
 498 variances for HP3O and PPO in the HalfCheetah environment for five different runs with different
 499 random seeds. HP3O has the highest explained variance over all runs suggesting that it is more robust
 500 against the variations of trajectories during learning. While for PPO, its explained variance can reach
 501 large negative values during training, which indicates the training instability when the trajectories
 502 vary significantly.

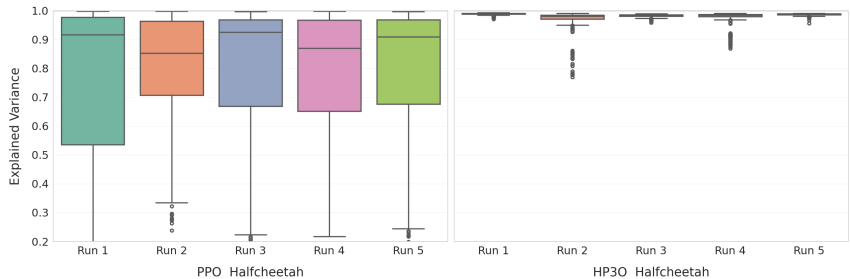


Figure 4: Explained Variance for HalfCheetah for PPO and HP3O

6.3 LIMITATIONS

514 Though theoretical and empirical results have shown that the proposed HP3O outperforms the
 515 popular baseline PPO over diverse control tasks, some limitations need to be discussed for potential
 516 improvement in the future. First, HP3O/HP3O+ require more hyperparameter tuning for the trajectory
 517 replay buffer, which can impact model performance compared to PPO. It has been acknowledged that
 518 hyperparameter tuning is critical for reinforcement learning such that for the hardest benchmarks, the
 519 already narrow basins of effective hyperparameters may become prohibitively small for our proposed
 520 algorithms, leading to poor performance. Second, in sparse reward environments, dequeuing the
 521 trajectory replay buffer can result in insufficient learning. Unlike the traditional replay buffer, which
 522 stores all experiences, our design requires the buffer to discard old trajectories so that the potential
 523 data distribution drift can be alleviated. This may cause a problem that good trajectories may only be
 524 learned once. Thus, the tradeoff between data distribution drift and learning frequency for the buffer
 525 needs to be investigated more in future work. Finally, there remains substantial room for performance
 526 improvement for the proposed algorithms compared to SAC. Further work in algorithm design is
 527 required to ensure HP3O/HP3O+ is on par with SAC but with low variance. The current ones can be
 528 regarded as one of the first steps toward bridging the gap between on-policy and off-policy methods.
 529

7 CONCLUSION AND BROADER IMPACTS

530 In this work, we presented a novel hybrid-policy reinforcement learning algorithm by incorporating a
 531 replay buffer into the popular PPO algorithm. Specifically, we utilized random sampling to reuse
 532 samples generated by the prior policies to improve the sample efficiency of PPO. We developed HP3O
 533 and theoretically derived its policy improvement lower bound. Subsequently, we designed a new
 534 advantage function in HP3O+ and presented a modified lower bound to provide theoretical guarantees.
 535 We investigated the stationary point convergence for HP3O and used several continuous control
 536 environments and baselines to showcase the superiority of the proposed algorithms. Additionally, we
 537 focused on variance reduction while maintaining high reward returns, encouraging the community to
 538 consider both high rewards and variance reduction. The theoretical claims of higher sample efficiency
 539 and variance reduction were empirically supported.

REFERENCES

- 540
541
542 Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In
543 *International conference on machine learning*, pp. 22–31. PMLR, 2017.
- 544 Fouad Bahrpeyma, Abishek Sunilkumar, and Dirk Reichelt. Application of reinforcement learning to
545 ur10 positioning for prioritized multi-step inspection in nvidia omniverse. In *2023 IEEE Symposium*
546 *on Industrial Electronics & Applications (ISIEA)*, pp. 1–6. IEEE, 2023.
- 547
548 Marco Biemann, Fabian Scheller, Xiufeng Liu, and Lizhen Huang. Experimental evaluation of
549 model-free reinforcement learning algorithms for continuous hvac control. *Applied Energy*, 298:
550 117164, 2021.
- 551 Johan Bjorck, Carla P Gomes, and Kilian Q Weinberger. Is high variance unavoidable in rl? a case
552 study in continuous control. *arXiv preprint arXiv:2110.11222*, 2021.
- 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 Jonathan Daniel Chang, Kianté Brantley, Rajkumar Ramamurthy, Dipendra Misra, and Wen Sun.
558 Learning to generate better than your large language models. 2023.
- 559 Xing Chen, Dongcui Diao, Hechang Chen, Hengshuai Yao, Haiyin Piao, Zhixiao Sun, Zhiwei Yang,
560 Randy Goebel, Bei Jiang, and Yi Chang. The sufficiency of off-policy and soft clipping: Ppo
561 is still insufficient according to an off-policy measure. In *Proceedings of the AAAI Conference on*
562 *Artificial Intelligence*, volume 37, pp. 7078–7086, 2023.
- 563
564 Jim G Dai and Mark Gluzman. Refined policy improvement bounds for mdps. *arXiv preprint*
565 *arXiv:2107.08068*, 2021.
- 566
567 Yan Dai, Haipeng Luo, Chen-Yu Wei, and Julian Zimmert. Refined regret for adversarial mdps with
568 linear function approximation. In *International Conference on Machine Learning*, pp. 6726–6759.
PMLR, 2023.
- 569
570 Stephen Dankwa and Wenfeng Zheng. Twin-delayed ddpq: A deep reinforcement learning tech-
571 nique to model a continuous movement of an intelligent robot agent. In *Proceedings of the 3rd*
572 *international conference on vision, image and signal processing*, pp. 1–5, 2019.
- 573
574 Debaprasad Dutta and Simant R Upreti. A survey and comparative evaluation of actor-critic methods
in process control. *The Canadian Journal of Chemical Engineering*, 100(9):2028–2056, 2022.
- 575
576 Rasool Fakoor, Pratik Chaudhari, and Alexander J Smola. P3o: Policy-on policy-off policy optimiza-
577 tion. In *Uncertainty in Artificial Intelligence*, pp. 1017–1027. PMLR, 2020.
- 578
579 Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-
critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- 580
581 Shixiang Shane Gu, Timothy Lillicrap, Richard E Turner, Zoubin Ghahramani, Bernhard Schölkopf,
582 and Sergey Levine. Interpolated policy gradient: Merging on-policy and off-policy gradient
583 estimation for deep reinforcement learning. *Advances in neural information processing systems*,
584 30, 2017.
- 585
586 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
587 maximum entropy deep reinforcement learning with a stochastic actor. In *International conference*
on machine learning, pp. 1861–1870. PMLR, 2018.
- 588
589 Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan,
590 John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning from demonstrations. In
591 *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 592
593 Yuenan Hou, Lifeng Liu, Qing Wei, Xudong Xu, and Chunlin Chen. A novel ddpq method with
prioritized experience replay. In *2017 IEEE international conference on systems, man, and*
cybernetics (SMC), pp. 316–321. IEEE, 2017.

- 594 Kai-Chun Hu, Chen-Huan Pi, Ting Han Wei, I Wu, Stone Cheng, Yi-Wei Dai, Wei-Yuan Ye,
595 et al. Towards combining on-off-policy methods for real-world applications. *arXiv preprint*
596 *arXiv:1904.10642*, 2019.
- 597
- 598 Mingzhe Hu, Jiahao Zhang, Luke Matkovic, Tian Liu, and Xiaofeng Yang. Reinforcement learning
599 in medical image analysis: Concepts, applications, challenges, and future directions. *Journal of*
600 *Applied Clinical Medical Physics*, 24(2):e13898, 2023.
- 601
- 602 Nai-Chieh Huang, Ping-Chun Hsieh, Kuo-Hao Ho, Hsuan-Yu Yao, Kai-Chun Hu, Liang-Chun
603 Ouyang, I Wu, et al. Neural ppo-clip attains global optimality: A hinge loss perspective. *arXiv*
604 *preprint arXiv:2110.13799*, 2021.
- 605
- 606 Ruinan Jin, Shuai Li, and Baoxiang Wang. On stationary point convergence of ppo-clip. In *The*
607 *Twelfth International Conference on Learning Representations*, 2023.
- 608
- 609 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance
610 reduction. *Advances in neural information processing systems*, 26, 2013.
- 611
- 612 Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In
613 *Proceedings of the Nineteenth International Conference on Machine Learning*, pp. 267–274, 2002.
- 614
- 615 Nathan Kallus and Masatoshi Uehara. Statistically efficient off-policy policy gradients. In *Interna-*
616 *tional Conference on Machine Learning*, pp. 5089–5100. PMLR, 2020.
- 617
- 618 Amandeep Kaur and Krishan Kumar. Energy-efficient resource allocation in cognitive radio networks
619 under cooperative multi-agent model-free reinforcement learning schemes. *IEEE Transactions on*
620 *Network and Service Management*, 17(3):1337–1348, 2020.
- 621
- 622 David M LaHuis, Michael J Hartman, Shotaro Hakoyama, and Patrick C Clark. Explained variance
623 measures for multilevel models. *Organizational Research Methods*, 17(4):433–451, 2014.
- 624
- 625 Shreyas Lele, Kavita Gangar, Harshal Daftary, and Dewashish Dharkar. Stock market trading agent
626 using on-policy reinforcement learning algorithms. *Available at SSRN 3582014*, 2020.
- 627
- 628 Timothée Lesort, Massimo Caccia, and Irina Rish. Understanding continual learning settings with
629 data distribution drift analysis. *arXiv preprint arXiv:2104.01678*, 2021.
- 630
- 631 Xingxing Liang, Yang Ma, Yanghe Feng, and Zhong Liu. Ptr-ppo: Proximal policy optimization with
632 prioritized trajectory replay. *arXiv preprint arXiv:2112.03798*, 2021.
- 633
- 634 Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
635 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv*
636 *preprint arXiv:1509.02971*, 2015.
- 637
- 638 Yanli Liu, Kaiqing Zhang, Tamer Basar, and Wotao Yin. An improved analysis of (variance-reduced)
639 policy gradient and natural policy gradient methods. *Advances in Neural Information Processing*
640 *Systems*, 33:7624–7636, 2020.
- 641
- 642 Yongshuai Liu, Avishai Halev, and Xin Liu. Policy learning with constraints in model-free reinforce-
643 ment learning: A survey. In *The 30th international joint conference on artificial intelligence (ijcai)*,
644 2021.
- 645
- 646 Daoming Lyu, Qi Qi, Mohammad Ghavamzadeh, Hengshuai Yao, Tianbao Yang, and Bo Liu.
647 Variance-reduced off-policy memory-efficient policy search. *arXiv preprint arXiv:2009.06548*,
2020.
- 648
- 649 Jincheng Mei, Wesley Chung, Valentin Thomas, Bo Dai, Csaba Szepesvari, and Dale Schuurmans.
650 The role of baselines in policy gradient optimization. *Advances in Neural Information Processing*
651 *Systems*, 35:17818–17830, 2022.
- 652
- 653 Wenjia Meng, Qian Zheng, Gang Pan, and Yilong Yin. Off-policy proximal policy optimization. In
654 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9162–9170, 2023.

- 648 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
649 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
650 *arXiv:1312.5602*, 2013.
- 651 Anh Tuan Nguyen, Duy Hoang Pham, Bee-Lan Oo, Mattheos Santamouris, Yonghan Ahn, and
652 Benson TH Lim. Modelling building hvac control strategies using a deep reinforcement learning
653 approach. *Energy and Buildings*, pp. 114065, 2024.
- 654 Matteo Papini, Damiano Binaghi, Giuseppe Canonaco, Matteo Pirota, and Marcello Restelli. Stochas-
655 tic variance-reduced policy gradient. In *International conference on machine learning*, pp. 4026–
656 4035. PMLR, 2018.
- 657 Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, Yun-Nung Chen, and Kam-Fai Wong. Adver-
658 sarial advantage actor-critic model for task-completion dialogue policy learning. In *2018 IEEE*
659 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6149–6153.
660 IEEE, 2018.
- 661 Athanasios S Polydoros and Lazaros Nalpantidis. Survey of model-based reinforcement learning:
662 Applications on robotics. *Journal of Intelligent & Robotic Systems*, 86(2):153–173, 2017.
- 663 Rafael Figueiredo Prudencio, Marcos ROA Maximo, and Esther Luna Colombini. A survey on offline
664 reinforcement learning: Taxonomy, review, and open problems. *IEEE Transactions on Neural*
665 *Networks and Learning Systems*, 2023.
- 666 James Queeney, Yannis Paschalidis, and Christos G Cassandras. Generalized proximal policy
667 optimization with sample reuse. *Advances in Neural Information Processing Systems*, 34:11909–
668 11919, 2021.
- 669 Dhruv Mauria Saxena, Sangjae Bae, Alireza Nakhaei, Kikuo Fujimura, and Maxim Likhachev.
670 Driving in dense traffic with model-free reinforcement learning. In *2020 IEEE International*
671 *Conference on Robotics and Automation (ICRA)*, pp. 5385–5392. IEEE, 2020.
- 672 Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv*
673 *preprint arXiv:1511.05952*, 2015.
- 674 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
675 policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR,
676 2015.
- 677 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
678 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 679 Uri Sherman, Tomer Koren, and Yishay Mansour. Improved regret for efficient online reinforcement
680 learning with linear function approximation. In *International Conference on Machine Learning*,
681 pp. 31117–31150. PMLR, 2023.
- 682 Marin Toromanoff, Emilie Wirbel, and Fabien Moutarde. End-to-end model-free reinforcement
683 learning for urban driving using implicit affordances. In *Proceedings of the IEEE/CVF conference*
684 *on computer vision and pattern recognition*, pp. 7153–7162, 2020.
- 685 Samuele Tosatto, Joao Carvalho, Hany Abdulsamad, and Jan Peters. A nonparametric off-policy
686 policy gradient. In *International Conference on Artificial Intelligence and Statistics*, pp. 167–177.
687 PMLR, 2020.
- 688 Chong Wang, Xi Chen, Alexander J Smola, and Eric P Xing. Variance reduction for stochastic
689 gradient optimization. *Advances in neural information processing systems*, 26, 2013.
- 690 Cathy Wu, Aravind Rajeswaran, Yan Duan, Vikash Kumar, Alexandre M Bayen, Sham Kakade,
691 Igor Mordatch, and Pieter Abbeel. Variance reduction for policy gradient with action-dependent
692 factorized baselines. *arXiv preprint arXiv:1803.07246*, 2018.
- 693 Pan Xu, Felicia Gao, and Quanquan Gu. Sample efficient policy gradient methods with recursive
694 variance reduction. *arXiv preprint arXiv:1909.08610*, 2019.

702 Pan Xu, Felicia Gao, and Quanquan Gu. An improved convergence analysis of stochastic variance-
703 reduced policy gradient. In *Uncertainty in Artificial Intelligence*, pp. 541–551. PMLR, 2020.

704 Andrea Zanette. When is realizability sufficient for off-policy reinforcement learning? In *International
705 Conference on Machine Learning*, pp. 40637–40668. PMLR, 2023.

706 Andrea Zanette, Ching-An Cheng, and Alekh Agarwal. Cautiously optimistic policy optimization and
707 exploration with linear function approximation. In *Conference on Learning Theory*, pp. 4473–4525.
708 PMLR, 2021.

709 Haijun Zhang, Ning Yang, Wei Huangfu, Keping Long, and Victor CM Leung. Power control
710 based on deep reinforcement learning for spectrum sharing. *IEEE Transactions on Wireless
711 Communications*, 19(6):4209–4219, 2020a.

712 Haijun Zhang, Minghui Jiang, Xiangnan Liu, Xiangming Wen, Ning Wang, and Keping Long.
713 Ppo-based pdacb traffic control scheme for massive iov communications. *IEEE Transactions on
714 Intelligent Transportation Systems*, 24(1):1116–1125, 2022.

715 Hang Zhang, Weike Liu, and Qingbao Liu. Reinforcement online active learning ensemble for
716 drifting imbalanced data streams. *IEEE Transactions on Knowledge and Data Engineering*, 34(8):
717 3971–3983, 2020b.

718 Junyu Zhang, Chengzhuo Ni, Csaba Szepesvari, Mengdi Wang, et al. On the convergence and sample
719 efficiency of variance-reduced policy gradient method. *Advances in Neural Information Processing
720 Systems*, 34:2228–2240, 2021.

721 Han Zhong and Tong Zhang. A theoretical analysis of optimistic proximal policy optimization in
722 linear markov decision processes. *Advances in Neural Information Processing Systems*, 36, 2024.

723 A APPENDIX

724 In this section, we present additional analysis and experimental results as a supplement to the main
725 contents. To conveniently refer to the theoretical results, we repeat the statements for all lemmas and
726 theorems.

727 A.1 ADDITIONAL THEORETICAL ANALYSIS

728 **Lemma 4.** (*Lemma 6.1 in Kakade & Langford (2002)*) For any policies $\hat{\pi}$ and π , we have

$$729 J(\hat{\pi}) - J(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\hat{\pi}}} [\mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} [A^{\pi}(s, a)]] \quad (7)$$

730 Lemma 4 signifies the cumulative return difference between two policies, π and $\hat{\pi}$.

731 **Lemma 5.** Consider any two policies $\hat{\pi}$ and π . Then the total variation distance between the state
732 visitation distributions $d^{\hat{\pi}}$ and d^{π} is bounded by

$$733 \delta(d^{\pi}, d^{\hat{\pi}}) \leq \frac{\gamma}{1 - \gamma} \mathbb{E}_{s \sim d^{\hat{\pi}}} [\delta(\pi, \hat{\pi})(s)], \quad (8)$$

734 where $\delta(\pi, \hat{\pi})(s)$ is defined in Lemma 1.

735 The proof follows similarly from Achiam et al. (2017). Next we present the proof for Lemma 2.

736 **Lemma 2:** Consider a present policy π_k , and any reference policy π_r . We then have, for any future
737 policy π ,

$$738 J(\pi) - J(\pi_k) \geq \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d^{\pi_r}} \left[\frac{\pi(a|s)}{\pi_r(a|s)} A^{\pi_k}(s, a) \right] - \frac{2\gamma C_{\pi_k}^{\pi}}{(1 - \gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)], \quad (9)$$

739 where $C_{\pi_k}^{\pi}$ and $\delta(\pi, \pi_r)(s)$ are defined as in Lemma 1.

740 *Proof.* The proof is similar to the proof of Lemma 7 in Queeney et al. (2021). We start from the
741 equality in Lemma 4 by adding and subtracting the term

$$742 \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] \quad (10)$$

743 □

With this, we obtain the following relationship:

$$\begin{aligned}
J(\pi) - J(\pi_k) &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] \\
&+ \frac{1}{1-\gamma} (\mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] - \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]]) \\
&\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] \\
&- \frac{1}{1-\gamma} |\mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] - \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]]|
\end{aligned} \tag{11}$$

The last inequality follows from the Triangle inequality. Subsequently, we can bound the second term of the last inequality using Hölder's inequality:

$$\begin{aligned}
\frac{1}{1-\gamma} |\mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] - \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]]| \\
\leq \frac{1}{1-\gamma} \|d^\pi - d^{\pi_r}\|_1 \|\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]\|_\infty,
\end{aligned} \tag{12}$$

where d^π and d^{π_r} both signify the state visitation distributions. In light of the definition of total variation distance and Lemma 3, the following relationship can be obtained accordingly

$$\|d^\pi - d^{\pi_r}\|_1 = 2\delta(d^\pi, d^{\pi_r}) \leq \frac{2\gamma}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)]. \tag{13}$$

Also note that

$$\|\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]\|_\infty = \max |\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]| = C_{\pi_k}^\pi. \tag{14}$$

Hence, substituting Eq. 13 and Eq. 14 into Eq. 12 and combining Eq. 11 yields the following inequality:

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] - \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)]. \tag{15}$$

Finally, without loss of generality, we assume that the support of π is contained in the support of π_r for all states, which is true for common policy representations used in policy optimization. We can rewrite the first term on the right hand side of the last inequality as

$$\frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi_r}} \left[\frac{\pi(a|s)}{\pi_r(a|s)} A^{\pi_k}(s, a) \right], \tag{16}$$

which leads to the desirable results.

Theorem 1: Consider prior policies $|\mathcal{B}|$ randomly sampled from the replay buffer R with indices $i = 0, 1, \dots, |\mathcal{B}| - 1$. For any distribution $v = [v_1, v_2, \dots, v_{|\mathcal{B}|}]$ over the $|\mathcal{B}|$ prior policies, and any future policy π generated by HP3O in Algorithm 1, the following relationship holds true

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{i \sim v} [\mathbb{E}_{(s,a) \sim d^{\pi_i}} \left[\frac{\pi(a|s)}{\pi_i(a|s)} A^{\pi_k}(s, a) \right]] - \frac{\gamma C_{\pi_k}^\pi \epsilon}{(1-\gamma)^2}, \tag{17}$$

where $C_{\pi_k}^\pi$ is defined as in Lemma 1.

Proof. Based on the definition of total variation distance, we have that

$$\mathbb{E}_{s \sim d^{\pi_k}} [\delta(\pi, \pi_k)(s)] = \mathbb{E} \left[\frac{1}{2} \int_{a \in \mathcal{A}} |\pi(a|s) - \pi_k(a|s)| da \right]. \tag{18}$$

We still make the assumption that the support of π is contained in the support of π_k for all states, which is true for the common policy representations used in policy optimization. Then, by multiplying and dividing by $\pi_k(a|s)$, we can observe that

$$\begin{aligned}
\mathbb{E}_{s \sim d^{\pi_k}} [\delta(\pi, \pi_k)(s)] &= \mathbb{E} \left[\frac{1}{2} \int_{a \in \mathcal{A}} \pi_k(a|s) \left| \frac{\pi(a|s)}{\pi_k(a|s)} - 1 \right| da \right] \\
&= \frac{1}{2} \mathbb{E}_{(s,a) \sim d^{\pi_k}} \left[\left| \frac{\pi(a|s)}{\pi_k(a|s)} - 1 \right| \right] \leq \frac{\epsilon}{2}.
\end{aligned} \tag{19}$$

The last inequality follows from the setup of PPO. With prior policies $\pi_i, i = 0, 1, 2, \dots, |\mathcal{B}| - 1$, we assume that the support of π is contained in the support of π_i for all states, which is true for common policy representations used in policy optimization. Based on Lemma 2, we can obtain

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi_i}} \left[\frac{\pi(a|s)}{\pi_i(a|s)} A^{\pi_k}(s, a) \right] - \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_i}} [\delta(\pi, \pi_i)(s)]. \tag{20}$$

Consider policy weights $v = [v_1, v_2, \dots, v_{|\mathcal{B}|}]$ over the policies in the minibatch \mathcal{B} . Thus, for any choice of distribution v , the convex combination determined by v of the $|\mathcal{B}|$ lower bounds given by

the last inequality yields the lower bound

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{i \sim v} [\mathbb{E}_{(s,a) \sim d^{\pi_i}} \left[\frac{\pi(a|s)}{\pi_i(a|s)} A^{\pi_k}(s, a) \right]] - \frac{2\gamma C^{\pi_k}}{(1-\gamma)^2} \mathbb{E}_{i \sim v} [\mathbb{E}_{s \sim d^{\pi_i}} [\delta(\pi, \pi_i)(s)]]]. \quad (21)$$

Combining Eq. 19 and Eq. 21, with some mathematical manipulation, results in the desirable conclusion. Now we're ready to prove Lemma 3. \square

Lemma 3: Consider a present policy π_k , and any reference policy π_r . We then have, for any future policy π ,

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi_r}} \left[\frac{\pi(a|s)}{\pi_r(a|s)} \hat{A}^{\pi_k}(s, a) \right] - \frac{2\gamma \hat{C}^{\pi_k}}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)] - \frac{2\gamma C^{\pi_k}}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)], \quad (22)$$

where $\hat{C}^{\pi_k} = \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} [\hat{A}^{\pi_k}(s, a)]|$, $\delta(\pi, \pi_r)(s)$ is defined as in Lemma 1, $C^{\pi_k} = \max_{s \in \mathcal{S}} |V^{\pi_k^*}(s) - V^{\pi_k}(s)|$.

Proof. Due to Lemma 1, we have

$$\begin{aligned} J(\pi) - J(\pi_k) &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [A^{\pi_k}(s, a)]] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [Q^{\pi_k}(s, a) - V^{\pi_k}(s)]] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [Q^{\pi_k}(s, a) - V^{\pi_k^*}(s) + V^{\pi_k^*}(s) - V^{\pi_k}(s)]]. \end{aligned} \quad (23)$$

Let $\hat{A}^{\pi_k}(s, a) = Q^{\pi_k}(s, a) - V^{\pi_k^*}(s)$ and $G^{\pi_k}(s) = V^{\pi_k^*}(s) - V^{\pi_k}(s)$ such that

$$J(\pi) - J(\pi_k) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [\hat{A}^{\pi_k}(s, a)]] + \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi(\cdot|s)} [G^{\pi_k}(s)]]. \quad (24)$$

Define $\|G^{\pi_k}(s)\|_\infty = \max_{s \in \mathcal{S}} |V^{\pi_k^*}(s) - V^{\pi_k}(s)| = C^{\pi_k}$. Follow similarly the proof from Lemma 2, we can attain the relationship as follows:

$$\begin{aligned} J(\pi) - J(\pi_k) &\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [\mathbb{E}_{a \sim \pi(\cdot|s)} [\hat{A}^{\pi_k}(s, a)]] + \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_r}} [G^{\pi_k}(s)] \\ &\quad - \frac{2\gamma \hat{C}^{\pi_k}}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)] \\ &\quad - \frac{2\gamma C^{\pi_k}}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_r}} [\delta(\pi, \pi_r)(s)]. \end{aligned} \quad (25)$$

The fact that $\min_{s \in \mathcal{S}} |V^{\pi_k^*}(s) - V^{\pi_k}(s)| = 0$ retains the desirable result.

Theorem 2: Consider prior policies $|\mathcal{B}|$ randomly sampled from the replay buffer R with indices $i = 0, 1, \dots, |\mathcal{B}| - 1$. For any distribution $v = [v_1, v_2, \dots, v_{|\mathcal{B}|}]$ over the $|\mathcal{B}|$ prior policies, and any future policy π generated by HP3O+ in Algorithm 1, the following relationship holds true

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{i \sim v} [\mathbb{E}_{(s,a) \sim d^{\pi_i}} \left[\frac{\pi(a|s)}{\pi_i(a|s)} \hat{A}^{\pi_k}(s, a) \right]] - \frac{\gamma \hat{C}^{\pi_k} \epsilon}{(1-\gamma)^2} - \frac{\gamma C^{\pi_k} \epsilon}{(1-\gamma)^2}, \quad (26)$$

where \hat{C}^{π_k} and C^{π_k} are defined as in Lemma 3.

Proof. Following the proof techniques in Theorem 1 and combining the conclusion from Lemma 3 obtains Eq. 26. \square

864 A.2 RISK OF OVERFITTING? 865

866 In our approach, each set of sampled trajectories includes the current best action trajectory in the
867 buffer, but we use a uniform distribution to sample mini-batch data points from all the trajectories
868 rather than only focusing on the best one. Additionally, the number of sampled trajectories is a
869 tunable parameter that we adjust based on the specific environment. Therefore, we ensure that the
870 model is exposed to a diverse set of experiences, which also helps mitigate the risk of overfitting.
871 Another important point is that our trajectory buffer operates on a FIFO (FirstIn-First-Out) basis.
872 As newer trajectories are added to the buffer, the oldest ones are replaced. This buffer maintains a
873 dynamic structure where trajectories are continually updated to reflect the most recent learning and
874 also helps to reduce distribution drift. We expect that these newer trajectories are more likely to be
875 better-performing as they are generated from the most current learned policy. All these techniques
876 are implemented in our buffer and help to balance exploration with prioritizing higher-performing
877 trajectories while also reducing the risk of overfitting.

878 A.3 INCORPORATION OF THE WORST TRAJECTORIES 879

880 In our approach, we prioritize leveraging higher-performing trajectories to optimize the agent’s
881 learning efficiency and to accelerate convergence toward optimal policies. This focus allows the agent
882 to reinforce successful behaviors more effectively. However, we understand the concern regarding
883 forgetting catastrophic behaviors, which could potentially lead to the agent’s catastrophic behaviors.
884 In practice, the FIFO buffer and uniform sampling from the sampled trajectories make sure that
885 a diverse range of experiences, including suboptimal or catastrophic behaviors, are preserved to
886 some extent within the buffer. This diversity helps the agent to maintain a broad understanding
887 of the environment, including both successful and unsuccessful strategies. Additionally, while we
888 do not explicitly prioritize the worst trajectories, our approach does not entirely discard them. By
889 maintaining a diverse buffer, the agent is still exposed to these behaviors, which can serve as alerting
890 examples. This exposure helps the agent learn to avoid repeating such catastrophic actions without
891 the need to focus on the worst trajectories explicitly. We believe this balance allows the agent to focus
892 on learning from successful strategies while still retaining an understanding of less optimal behaviors,
893 reducing the risk of catastrophic forgetting.

894 A.4 SAMPLE EFFICIENCY ANALYSIS 895

896 In this section, we present the sample efficiency analysis for the proposed HP3O algorithm, compared
897 to the vanilla PPO algorithm, which remains the most popular on-policy scheme so far. Though
898 the analysis is conducted particularly for the comparison between PPO and HP3O, the techniques
899 apply extensively to other on-policy policy-gradient-based algorithms whenever they satisfy the
900 conservative policy iteration property Kakade & Langford (2002); Achiam et al. (2017) to have
901 the policy improvement lower bounds. In this study, we aim to show how the off-policy sample
902 reuse significantly affects the original sample efficiency PPO has. We will not directly show the
903 exact sample complexity of HP3O and the improvement on top of PPO. For instance, to arrive at an
904 ϵ -optimality for policy gradient-based algorithms, a few works Zhong & Zhang (2024); Zanette et al.
905 (2021); Dai et al. (2023); Sherman et al. (2023) have revealed the exact complexity with respect to ϵ ,
906 but only for MDPs with linear function approximation. The exact sample complexity analysis for
907 the on-policy PPO algorithm remains extremely challenging and requires a substantial amount of
908 non-trivial effort. Thereby, in this paper, we disclose the impact of off-policy sample reuse on the
909 tradeoff between sample efficiency and learning stability.

910 To start with the comparison between PPO and HP3O, we denote by ϵ_H and ϵ_P the clipping parameters
911 for HP3O and PPO. Such a clipping parameter indicates the worst-case expected performance loss of
912 update at every time step. We next present a lemma that shows the relationship between ϵ_H and ϵ_P .

913 **Lemma 6.** *Consider prior policies $|\mathcal{B}|$ randomly sampled from the replay buffer R with indices
914 $i = 0, 1, \dots, |\mathcal{B}| - 1$. For any distribution $v = [v_1, v_2, \dots, v_{|\mathcal{B}|}]$ over the $|\mathcal{B}|$ prior policies, both HP3O
915 and PPO have the same worst-case expected performance loss at every update when the clipping
916 parameters satisfy the following condition:*

$$917 \epsilon_H = \frac{\epsilon_P}{\mathbb{E}_{i \sim v}[i + 1]}. \quad (27)$$

918 *Proof.* Recall from PPO such that

$$919 \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{s \sim d^{\pi_k}} [\delta(\pi, \pi_k)(s)] \leq \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \frac{\epsilon_P}{2}. \quad (28)$$

921 For HP3O, its penalty term in the policy improvement lower bound in Theorem 1 can be upper
922 bounded by using the Triangle inequality. Therefore, we have the following relationship

$$923 \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{i \sim v} [\mathbb{E}_{s \sim d^{\pi_i}} [\delta(\pi, \pi_i)(s)]] \\ 924 \leq \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{i \sim v} \left[\sum_{j=0}^i \mathbb{E}_{s \sim d^{\pi_i}} [\delta(\pi_{j+1}, \pi_j)(s)] \right]. \quad (29)$$

928 The last inequality holds if the prior policies are in a chronological order based on their histories.
929 In practice, we do not set such an order for them, but due to the FIFO strategy we have leveraged,
930 they can still be set in this for the sake of analysis. Since we still resort to the clipping mecha-
931 nism in HP3O, each policy update approximately bounds each expected total variation distance
932 $\mathbb{E}_{s \sim d^{\pi_i}} [\delta(\pi_{j+1}, \pi_j)(s)]$ by $\frac{\epsilon_H}{2}$, which follows analogously from that in PPO. With this in hand, we
933 are now able to further bound Eq. 29 in the following relationship

$$934 \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{i \sim v} [\mathbb{E}_{s \sim d^{\pi_i}} [\delta(\pi, \pi_i)(s)]] \\ 935 \leq \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \mathbb{E}_{i \sim v} \left[\frac{\epsilon_H}{2} (i+1) \right] \quad (30) \\ 936 \leq \frac{2\gamma C_{\pi_k}^\pi}{(1-\gamma)^2} \frac{\epsilon_H}{2} \mathbb{E}_{i \sim v} [i+1]$$

940 Comparing the bounds in Eq. 28 and Eq. 30 yields the desirable result. \square

942 Lemma 6 technically shows us that if the two clipping parameters ϵ_H and ϵ_P satisfy the condition of
943 $\epsilon_H = \frac{\epsilon_P}{\mathbb{E}_{i \sim v} [i+1]}$, the worst-case expected performance loss at each update remains roughly the same.
944 This intuitively makes sense as HP3O leverages prior policies from the replay buffer to update the
945 policy model, which requires it to perform smaller updates. A benefit from this is to make policy
946 updates more frequently, thus schematically stabilizing policy learning. In what follows, we present
947 more analysis about this tradeoff.

948 To ease the analysis, we assume that the policies in the training batch \mathcal{B} are randomly sampled with
949 uniform policy weights, i.e., $v_i = \frac{1}{|\mathcal{B}|}$, for $i = 0, 1, \dots, |\mathcal{B}| - 1$, for collecting data to train the network
950 models. However, more advanced techniques such as Prioritized Experience Replay (PER) Schaul
951 et al. (2015) can be applied accordingly. In each episode, we also assume that for PPO, it requires
952 $N = Mn$ samples for sufficiently training the critic and actor networks, where M is the number of
953 mini-batch and n is the batch size. In this setting, PPO makes one episodic update upon the current
954 policy π_k by traversing N samples generated by π_k . However, for HP3O, since there exist multiple
955 policies prior to π_k , it is able to make M updates sourced from different prior policies per N samples
956 collected from \mathcal{B} , as long as $|\mathcal{B}| \leq M$. Thus, we next show that HP3O is able to increase the change
957 in the total variational distance of the policy throughout training, without sacrificing stability, when it
958 is compared to PPO.

959 **Theorem 3.** Suppose that $|\mathcal{B}| = M$ and that the policies in the training batch \mathcal{B} are randomly
960 sampled with uniform policy weights, i.e., $v_i = \frac{1}{|\mathcal{B}|}$, for $i = 0, 1, \dots, |\mathcal{B}| - 1$. Then, HP3O has a
961 larger frequency of change in total variation distance of the policy throughout training by a factor of
962 $\frac{2M}{M+1}$ compared to PPO, while using the same number of samples for each update as PPO.

964 *Proof.* Pertaining to Lemma 6 and the fact that $|\mathcal{B}| = M$, we have the following relationship:

$$965 \epsilon_H = \frac{\epsilon_P}{\frac{1}{M} \sum_{i=0}^{M-1} (i+1)} = \frac{2\epsilon_P}{M+1}. \quad (31)$$

967 PPO makes one episodic policy update after N samples are collected, say from k to $k+1$, which
968 yields a policy change of $\frac{\epsilon_P}{2}$ in terms of the total variation distance. While for HP3O, it resorts to
969 data from prior policies to obtain N samples and makes M policy updates, as mentioned before. This
970 results in the overall policy change of

$$971 M \frac{\epsilon_H}{2} = \frac{2M}{M+1} \frac{\epsilon_P}{2}. \quad (32)$$

Thus, HP3O has a larger frequency of changes in the total variation distance of the policy throughout training by a factor of $\frac{2M}{M+1}$ compared to PPO, with the same number of samples. \square

By far, we have discussed the tradeoff between learning stability and sample size that biases toward learning stability when maintaining the same sample size as in PPO. Alternatively, we can perceive the problem from another perspective, in which HP3O needs to increase the sample size while maintaining the same change in total variation distance throughout training. A formal result is summarized as follows.

Theorem 4. *Suppose that $|\mathcal{B}| = 2M - 1$ and that the policies in the training batch \mathcal{B} are randomly sampled with uniform policy weights, i.e., $v_i = \frac{1}{|\mathcal{B}|}$, for $i = 0, 1, \dots, |\mathcal{B}| - 1$. Thus, HP3O increases the sample size used for each policy update by a factor of $\frac{2M-1}{M}$ compared to PPO, simultaneously maintaining the same change in the total variation of distance of the policy throughout training as PPO.*

Proof. As $|\mathcal{B}| = 2M - 1$, HP3O uses $(2M - 1)n$ to calculate each policy update from the prior to the new policy, compared to Mn samples used in PPO. Hence, HP3O increases the sample size used for each policy update by a factor of $\frac{2M-1}{M}$ compared to PPO. Immediately, based on Lemma 6, we can obtain

$$\epsilon_H = \frac{\epsilon_P}{\frac{1}{2M-1} \sum_{i=0}^{2M-2} (i+1)} = \frac{\epsilon_P}{M}. \quad (33)$$

We have shown in Theorem 3 that PPO makes one policy update with N samples collected, while HP3O makes M policy updates with the same number of samples collected. We then have

$$\frac{M\epsilon_H}{2} = \frac{\epsilon_P}{2}. \quad (34)$$

This implies that the overall change in total variation distance in HP3O is the same as in PPO. \square

One implication from Theorem 3 and 4 is that HP3O with uniform policy weights enhances the tradeoff between learning stability and sample efficiency in the vanilla PPO when $|\mathcal{B}|$ is selected between $[M, 2M - 1]$. This also motivates us to set the FIFO strategy as the selected training batch \mathcal{B} cannot deviate too far away from the current policy. Otherwise, the negative impact of distribution drift could be extreme.

A.5 HP3O vs. HP3O+

In the last subsection, we have shown that HP3O enables more frequent changes in the total variational distance of the policy throughout training, with the smaller updates. Though more changes in the total variational distance of the policy may help improve the sample efficiency, but in order to address the distribution drift, smaller updates are the resulting outcome, possibly slowing down the convergence. Hence, introducing the best trajectory π^* in HP3O+ assists in mitigating this issue. Since it can increase the update, while maintaining the same number of changes as in HP3O. Such a behavior is empirically shown to enhance the model performance. We summarize the larger update in the total variational distance in a formal theoretical result as follows.

Theorem 5. *Denote by D_{TV}^H and D_{TV}^{H+} the updates of total variational distance of the policies for HP3O and HP3O+, respectively, at the time step k . Then we have $D_{TV}^H \leq D_{TV}^{H+}$ for all k .*

Proof. In light of Theorem 1 and Theorem 2, we know that $D_{TV}^H = \frac{\gamma C_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2}$ and $D_{TV}^{H+} = \frac{\gamma \hat{C}_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2} + \frac{\gamma C_{\pi_k}^{\pi} \epsilon}{(1-\gamma)^2}$. We next show the latter is bounded below by the former. As $A^{\pi_k}(s, a) = \hat{A}^{\pi_k}(s, a) + G^{\pi_k}(s)$, $\hat{A}^{\pi_k}(s, a) = Q^{\pi_k}(s, a) - V^{\pi_k^*}(s)$, and $G^{\pi_k}(s) = V^{\pi_k^*}(s) - V^{\pi_k}(s)$, we have the following relationship

$$A^{\pi_k}(s, a) = \hat{A}^{\pi_k}(s, a) + V^{\pi_k^*}(s) - V^{\pi_k}(s) \quad (35)$$

Taking the expectation of the action $a \sim \pi(\cdot|s)$ on both sides yields

$$\mathbb{E}_{a \sim \pi(\cdot|s)}[A^{\pi_k}(s, a)] = \mathbb{E}_{a \sim \pi(\cdot|s)}[\hat{A}^{\pi_k}(s, a)] + V^{\pi_k^*}(s) - V^{\pi_k}(s), \quad (36)$$

1026 which leads to

$$1027 \quad |\mathbb{E}_{a \sim \pi(\cdot|s)}[A^{\pi_k}(s, a)]| = |\mathbb{E}_{a \sim \pi(\cdot|s)}[\hat{A}^{\pi_k}(s, a)] + V^{\pi_k^*}(s) - V^{\pi_k}(s)|$$

$$1028 \quad \leq |\mathbb{E}_{a \sim \pi(\cdot|s)}[\hat{A}^{\pi_k}(s, a)]| + |V^{\pi_k^*}(s) - V^{\pi_k}(s)|. \quad (37)$$

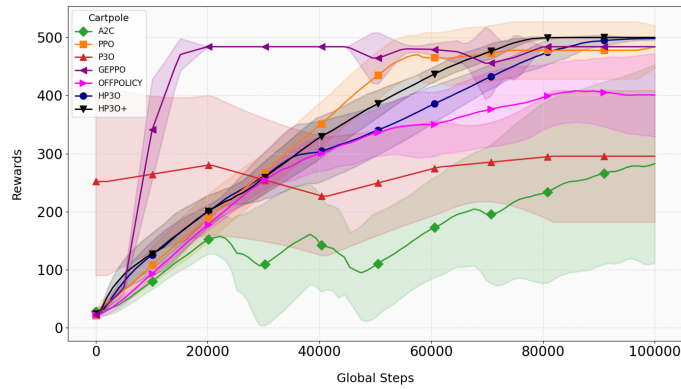
1029 This last inequality is due to the Triangle inequality. Taking the maximum operator of the state $s \sim \mathcal{S}$
 1030 on both sides results in the following:

$$1031 \quad \max_{s \sim \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)}[A^{\pi_k}(s, a)]| \leq \max_{s \sim \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)}[\hat{A}^{\pi_k}(s, a)]| + \max_{s \sim \mathcal{S}} |V^{\pi_k^*}(s) - V^{\pi_k}(s)|. \quad (38)$$

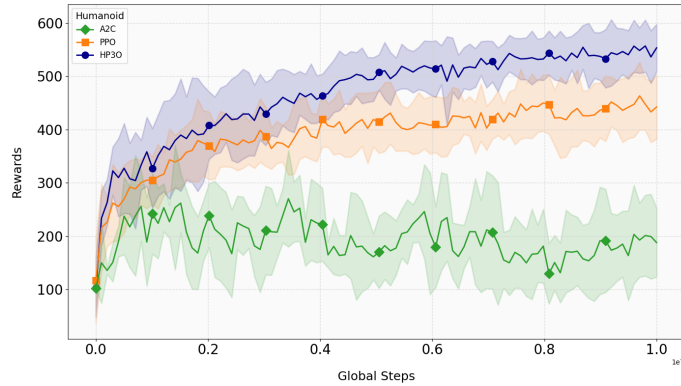
1032 Multiplying both sides in the above inequality by $\frac{\gamma\epsilon}{(1-\gamma)^2}$ yields the desirable result. \square

1033 A.6 TRAINING RESULTS FOR OTHER ENVIRONMENTS

1034 The following plot in Figure 5a presents the training curves obtained by training both the baseline
 1035 algorithms and our policy. These results further support our claim in the main paper that our policy
 1036 reduces variance while maintaining a high reward at the end.



1042 (a) Training performance of HP30 and PPO on the Cartpole environment
 1043 over 100k steps.



1044 (b) Training performance of HP30 and PPO on the Humanoid environ-
 1045 ment over 10 million steps.

1046 Figure 5: Comparison of HP30 and PPO training curves across different environments. (a) shows the
 1047 performance on Cartpole, while (b) shows the performance on Humanoid.

1048 A.7 ADDITIONAL EXPERIMENTAL RESULTS

1049 **Definition of explained variance.** The explained variance (EV) measures the proportion to which
 a mathematical model accounts for the variation of a given data set, which can be mathematically

defined in the following:

$$EV = 1 - \frac{Var(y - \hat{y})}{Var(y)}, \quad (39)$$

where y is the groundtruth and \hat{y} is the prediction. EV values typically vary from 0 to 1. In some scenarios, the value may be a large negative number, which indicates a poor prediction of y . Explained variance is a well-known metric in reinforcement learning, particularly for assessing the accuracy of value function predictions. In our experiment, explained variance was used to evaluate how well the value function predicts actual returns. The different runs correspond to separate training instances with different random seeds. The explained variance score is a risk metric that measures the dispersion of errors in a dataset. A score closer to 1.0 is better, as it indicates smaller squares of standard deviations of errors.

A.8 EXPLAINED VARIANCE FOR OTHER ENVIRONMENTS

Explained variance is a well-known metric in reinforcement learning, particularly for assessing the accuracy of value function predictions. In our experiment, explained variance was used to evaluate how well the value function predicts actual returns. The different runs correspond to separate training instances with different random seeds.

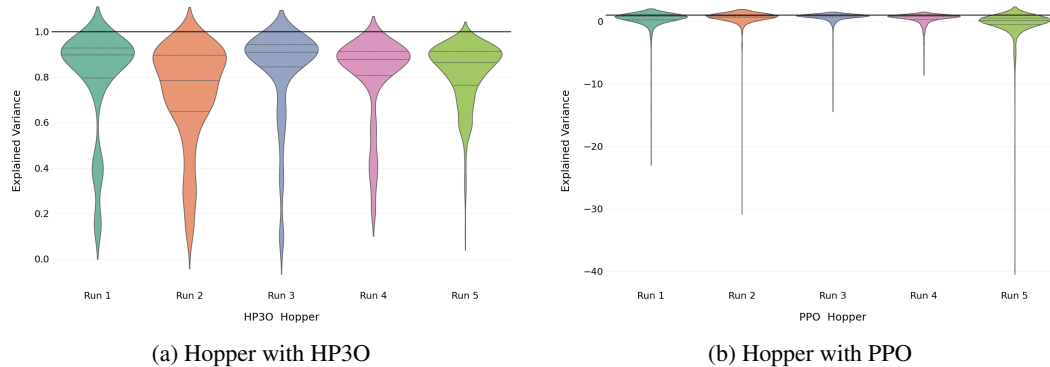


Figure 6: Explained Variance for Hopper

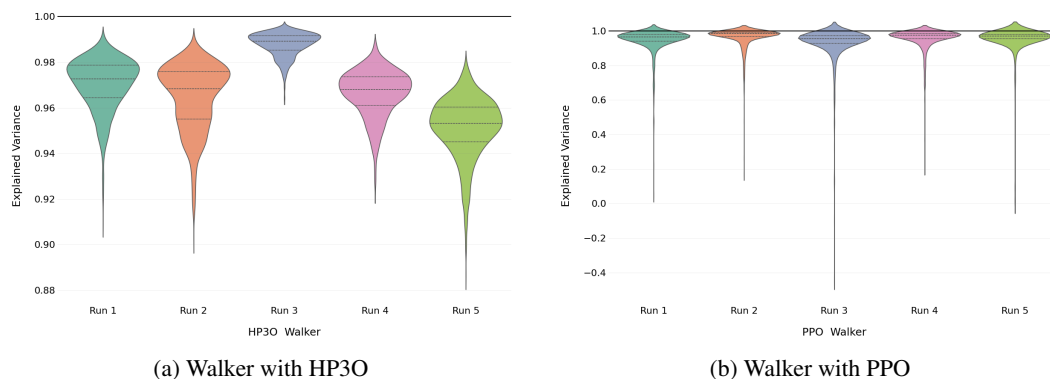


Figure 7: Explained Variance for Walker

A.9 SAC TRAINING BENCHMARKS

The following plots showcase the benchmark training results obtained by using the SAC policy. In some environments, SAC shows a relatively large variance. A notable disadvantage of SAC is that it only works with continuous action spaces.

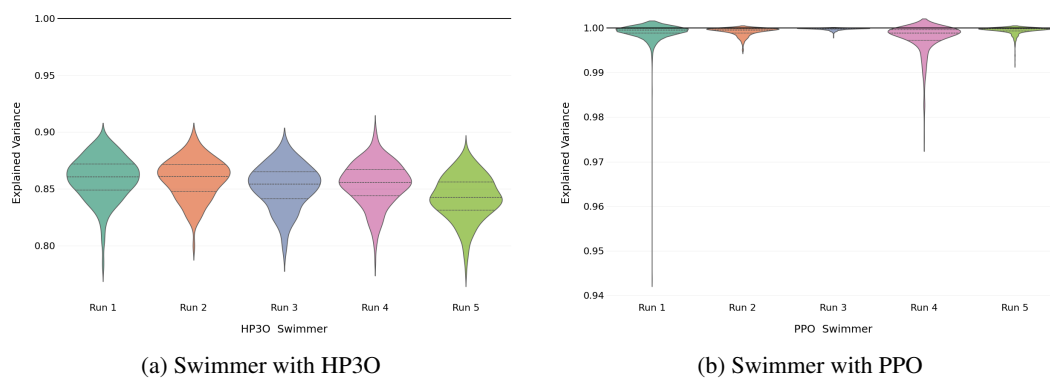


Figure 8: Explained Variance for Swimmer

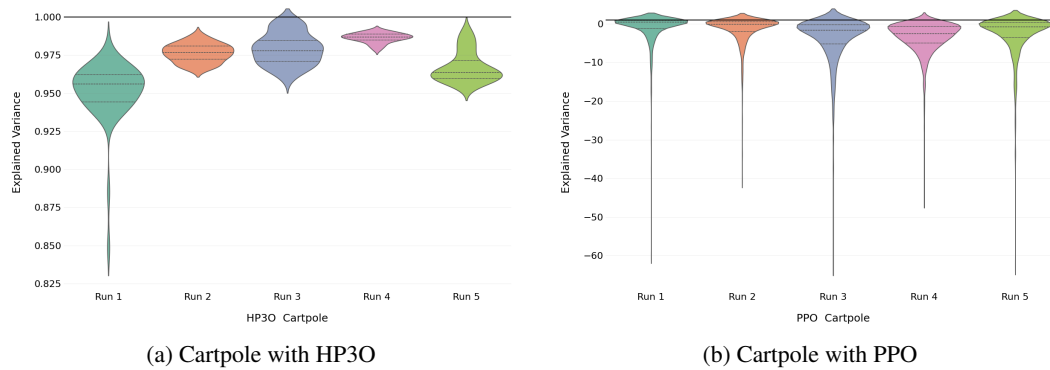


Figure 9: Explained Variance for Cartpole

A.10 EXPERIMENTAL CONFIGURATION

Experiments were performed on a local machine with an Intel Core i7-14700 CPU, 128 GB of RAM, and NVIDIA RTX 4090 GPU. We provide detailed information on our algorithms’ hyperparameters for all environments in our GitHub repository, which will be released once the paper is published.

A.11 DATA STORAGE FOR RL TRAINING

- **Trajectory Buffer:** Stores complete trajectories τ as sequences of (s_t, a_t, r_t, s_{t+1}) .
 - **Data Types:** Arrays of states, actions, rewards, and next states.
 - **Dimensions:**
 - * States s_t : Typically \mathbb{R}^n where n is the dimension of the state space.
 - * Actions a_t : Depends on the action space, usually \mathbb{R}^m where m is the dimension of the action space.
 - * Rewards r_t : Scalar values.
 - * Next states s_{t+1} : Same as states s_t .

A.12 CONVERGED REWARD

In this subsection, we present the converged reward at the last time step. In some environments, the training curves do not fully converge, which may make it challenging to assess the ultimate performance. However, for consistency across all algorithms, we maintained the same number of training time steps for each experiment. This allows for a fair comparison of sample efficiency across different methods, even if the algorithms did not always fully converge within the given time frame.

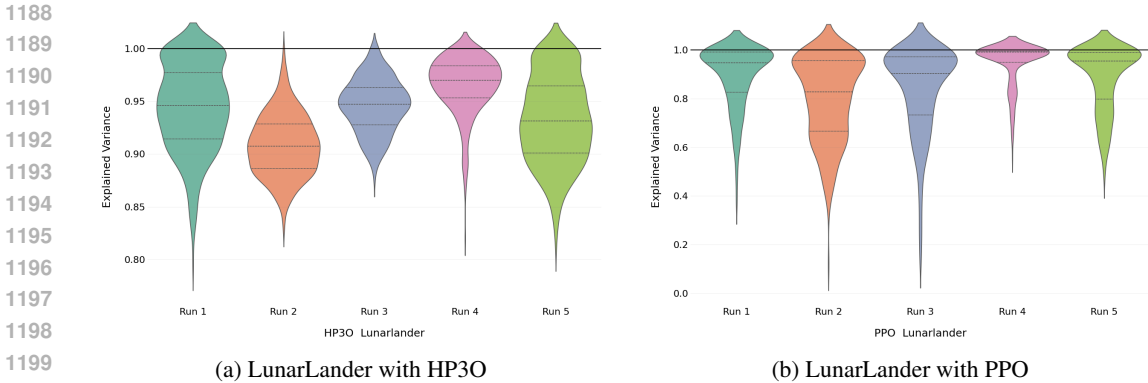


Figure 10: Explained Variance for LunarLander

Additionally, in some environments, we do present converged training curves, demonstrating the capabilities of the algorithms. In the reinforcement learning community, it is a common practice to show learning curves at a fixed number of steps for comparative analysis, even if full convergence is not always achieved. Notably, papers on SAC Haarnoja et al. (2018), PPO Schulman et al. (2017), GEppo Queeney et al. (2021), and Off-Policy PPO Meng et al. (2023) follow similar practices, with many of the environments presented in these works employing non-converged curves to provide valuable insights into training dynamics and sample efficiency. Table 2 shows the detailed converged reward performance of all different RL algorithms over different continuous tasks. In order to ensure a fair comparison between GEppo and our method, we first analyzed the performance differences between the PPO baselines in our implementation vs. in the GEppo repository. These discrepancies were primarily due to variations in implementation details (e.g., leveraging TensorFlow packages, early version of Mujoco environment), which significantly impacted the baseline performance. To address the discrepancies, we normalized the PPO baseline results to match our implementation. As a result, both baselines produced comparable outcomes. We then applied the same normalization factor to the GEppo results, repeating this procedure for each environment to ensure fair and consistent comparisons across all settings.

Table 2: Summary of mean and standard deviation of rewards for each policy across diverse environments (in the form Mean \pm Std) at or close to the converged stage. The bold one represents the best reward performance.

Environment	A2C	GEppo	HP3O	HP3O+	OffPolicy	P3O	PPO
CartPole	282.47 \pm 170.87	21.76 \pm 2.54	498.25 \pm 2.51	500.00 \pm 0.00	400.31 \pm 72.60	295.38 \pm 113.44	483.59 \pm 36.70
HalfCheetah	334.16 \pm 302.14	2156.31 \pm 1024.06	3523.20 \pm 565.39	3967.47 \pm 244.80	738.11 \pm 274.89	1251.17 \pm 517.84	2276.87 \pm 902.20
Hopper	120.31 \pm 48.29	976.33 \pm 68.36	988.67 \pm 16.53	1891.35 \pm 79.47	961.64 \pm 76.69	1107.49 \pm 281.71	946.90 \pm 64.86
InvertedPendulum	71.26 \pm 76.20	463.02 \pm 0.00	956.66 \pm 36.20	1000.00 \pm 0.00	11.24 \pm 3.22	810.77 \pm 46.26	463.02 \pm 445.28
LunarLander	-69.13 \pm 105.52	136.07 \pm 12.33	225.91 \pm 10.97	146.63 \pm 7.10	114.68 \pm 106.51	-624.72 \pm 292.51	130.58 \pm 16.53
Swimmer	12.99 \pm 6.25	133.83 \pm 6.69	340.00 \pm 4.18	343.40 \pm 1.41	157.73 \pm 107.07	32.75 \pm 14.32	131.98 \pm 38.76
Walker	134.05 \pm 51.54	1140.31 \pm 389.75	1934.91 \pm 152.69	1895.29 \pm 98.74	2093.24 \pm 370.08	1777.91 \pm 465.01	1150.36 \pm 97.18

A.13 COMPUTATIONAL EFFICIENCY

To probe particularly the computational efficiency of diverse algorithms presented in this study, we compare them in the wall-clock time spent to reach a certain reward in the HalfCheetah environment. Due to the time limitation, we are unable to obtain results for all other environments, while including them in the final version. Such an investigation offers us useful insights about which methods are more practically feasible and deployable given the limited real-time budget. Figure 12 shows the specific performance of wall-clock time cost for different approaches reaching the rewards of 1100 and 2100, respectively. One immediate observation from the results is that GEppo requires significantly more time to converge compared to all other schemes. For a couple of algorithms, such as A2C and OffPolicy, the training progresses are pretty slow, eventually failing to achieve rewards of 1100 and 2100 in the HalfCheetah environment. Another implication of interest from the results is that at the beginning, PPO may progress faster, compared to both HP3O and HP3O+. However, due to its on-policy behavior, the sample inefficiency issue still affects the overall training progress. Different

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

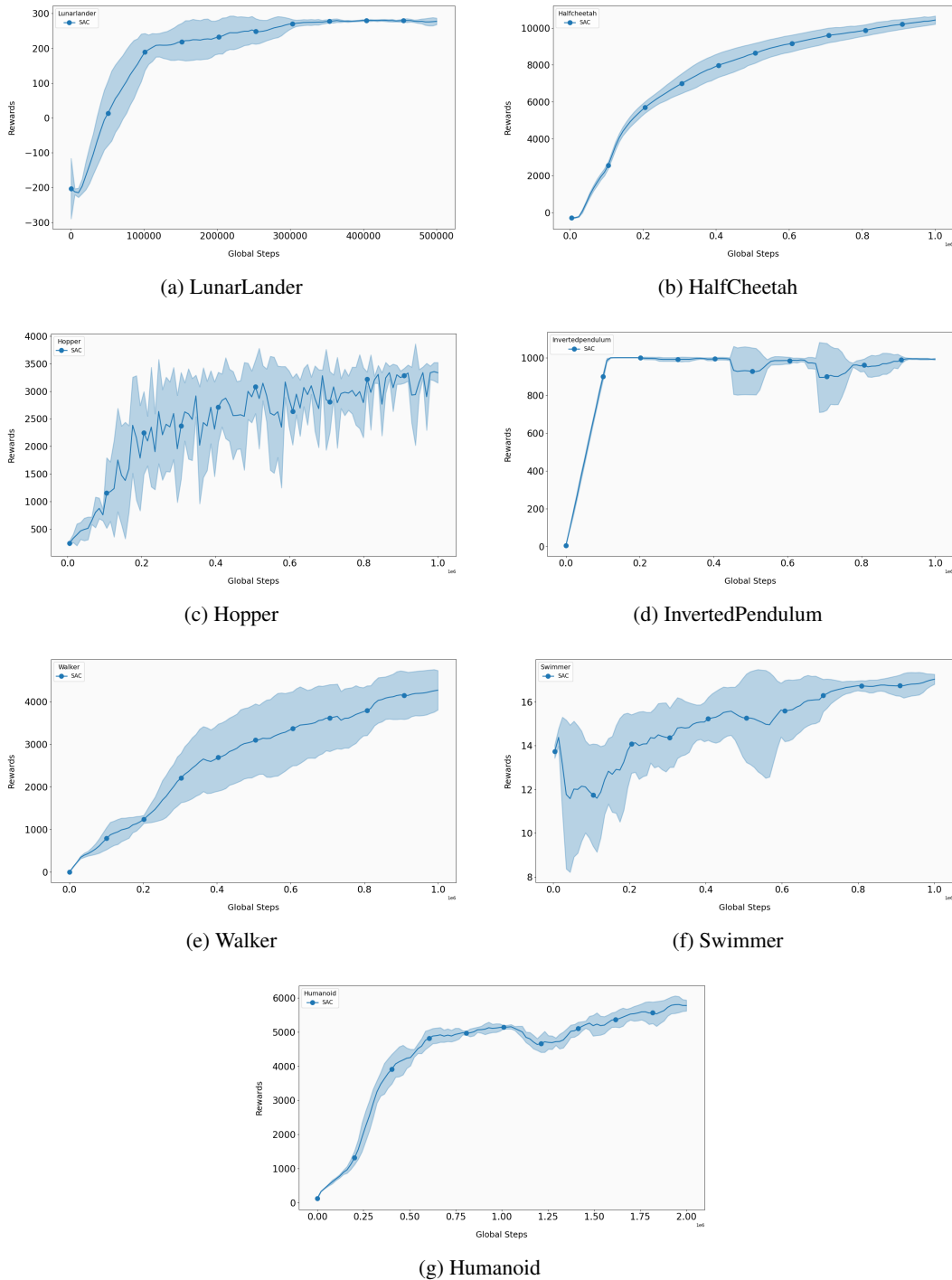
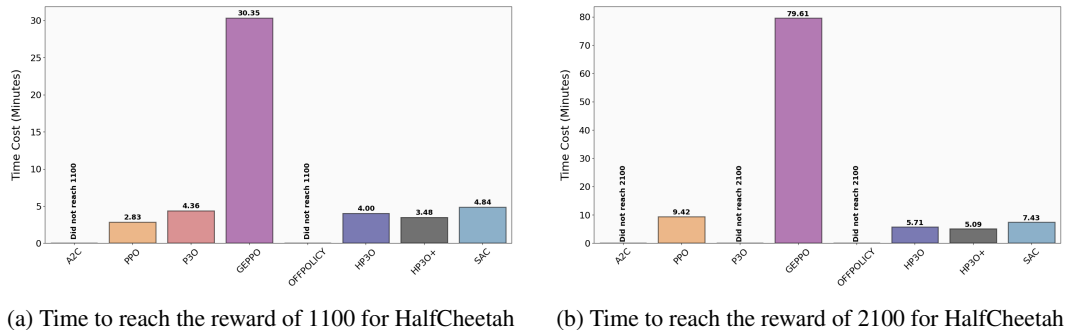


Figure 11: Training curves for SAC across various environments. The solid curves indicate the mean, while the shaded areas represent the standard deviation over the five runs.

from that, both HP3O and HP3O+ make consistent progress throughout the training process and take minimal time to achieve certain rewards. Between them, HP3O+ has slightly better performance, which empirically validates our conclusion from Theorem 5. The finding also complies with that in Figure 3b.

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



(a) Time to reach the reward of 1100 for HalfCheetah (b) Time to reach the reward of 2100 for HalfCheetah

Figure 12: Time taken by algorithms to reach specific reward thresholds in the HalfCheetah environment.

A.14 IMPACT OF BUFFER SIZE AND MEMORY USAGE

In this subsection, we also discuss the impact of the buffer size on the model performance qualitatively. The other two aspects included here are the discussion on the sampling strategy and the trade-off in sparse reward settings.

- **Large Trajectory Buffer Size:** A larger trajectory buffer size allows us to store a greater number of diverse trajectories. This diversity can enhance generalization and reduce overfitting, as the agent learns from a broad range of experiences. However, in sparse reward settings, maintaining a large trajectory buffer may mean that the inclusion of outdated or less relevant trajectories could introduce instability and slow down learning, as the agent may be exposed to experiences that no longer align with its current policy.
- **Small Trajectory Buffer Size:** A smaller trajectory buffer retains fewer trajectories, which typically results in the agent learning from more recent experiences that are closely aligned with the current policy. This can improve stability, as updates are based on recent, relevant data. However, a smaller buffer can reduce the diversity of sampled experiences, leading to an increased risk of overfitting and limiting the agent’s ability to effectively explore different parts of the environment.
- **Sampling Strategy:** Our sampling strategy also plays a critical role in managing the trade-off between stability and performance. By ensuring that the best return trajectory is always included in the sampled trajectories, we provide a strong guiding signal that improves policy performance. The sample rate, where we evenly sample from each selected trajectory, helps in maintaining a balance between exploration and exploitation, as well as in utilizing high-quality trajectories effectively.
- **Trade-off in Sparse Reward Settings:** In sparse reward environments, the need to maintain high-quality trajectories becomes even more crucial. A large trajectory buffer can help capture rare, valuable experiences, but the risk of dequeuing these valuable trajectories before they can contribute meaningfully to learning is higher. Ensuring that the best trajectory is always sampled helps mitigate this issue, but the buffer size still influences how effectively these rare rewards are retained and leveraged.

Regarding memory consumption, the total memory usage depends on the following:

1. **Number of Trajectories in the Buffer (Buffer Size):** We use a fixed-size replay buffer, which holds N trajectories.
2. **Average Length of Trajectories (T):** Since trajectories have varying lengths, the memory usage will depend on the average trajectory length.
3. **Dimensionality of State and Action Spaces:** Let d_s be the dimension of the state, and d_a be the dimension of the action. Each step in a trajectory stores both state and action information.
4. **Data type size:** Denote by n the specific data type size.

The memory consumption (M) can be roughly estimated as:

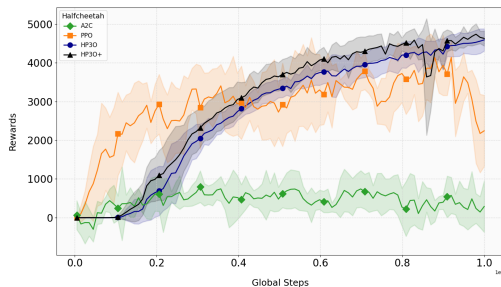
$$M = NTn(d_s + d_a)$$

For example, if $N = 1000$ trajectories are stored, each with an average length of 200 steps, and assuming $d_s = 17$, $d_a = 6$, and using 32-bit floats (4 bytes), the memory requirement would be:

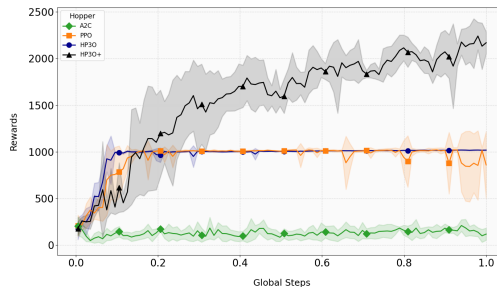
$$M = 1000 \times 200 \times (17 + 6) \times 4 \text{ bytes} = 18,400,000 \text{ bytes} \approx 18.4 \text{ MB}$$

This calculation provides an estimate, but the actual memory usage may vary depending on the environment and the specific implementation of the replay buffer.

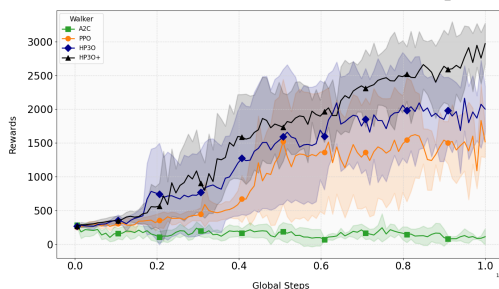
A.15 EVALUATION RESULTS



(a) Evaluation plot for the HalfCheetah environment.



(b) Evaluation plot for the Hopper environment.



(c) Evaluation plot for the Walker environment.

Figure 13: Evaluation plots for the HalfCheetah, Hopper, and Walker environments during the evaluation stage every 5000 steps. Each experiment includes five different runs with various random seeds. The solid curves indicate the mean, while the shaded areas represent the standard deviation over the five runs.

The evaluation results align with our training expectations. Overall, our presented models, HP3O and HP3O+, consistently outperform the baseline models across all environments, achieving higher rewards while maintaining relatively low variance. The PPO baseline performs well initially but tends to be less sample efficient and has relatively higher variance, whereas A2C struggles to reach comparable performance.

The results clearly demonstrate that our presented models, HP3O and HP3O+, are better equipped for these environments. This also verifies our claim from the training analysis. Both HP3O and HP3O+ combine improved sample efficiency and reduced variance, leading to more stable learning outcomes and higher cumulative rewards. These advantages enable our models to not only outperform the baselines but also maintain robustness and efficiency across diverse environments. Due to the time limitation, we are unable to obtain results for other environments with other methods, and will include additional results in the final version.