

GUIDED POLICY OPTIMIZATION UNDER PARTIAL OBSERVABILITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement Learning (RL) in partially observable environments poses significant challenges due to the complexity of learning under uncertainty. While additional information, such as that available in simulations, can enhance training, effectively leveraging it remains an open problem. To address this, we introduce Guided Policy Optimization (GPO), a framework that co-trains a guider and a learner. The guider takes advantage of privileged information while ensuring alignment with the learner’s policy that is primarily trained via imitation learning. We theoretically demonstrate that this learning scheme achieves optimality comparable to direct RL, thereby overcoming key limitations inherent in existing approaches. Empirical evaluations show strong performance of GPO across various tasks, including continuous control with partial observability and noise, and memory-based challenges, significantly outperforming existing methods.

1 INTRODUCTION

Many real-world tasks can be formulated as sequential decision-making problems where agents take actions in an environment to achieve specific goals over time (Puterman, 2014). Reinforcement Learning (RL) has emerged as a powerful tool for solving such tasks, leveraging trial-and-error learning to optimize long-term rewards (Sutton & Barto, 2018). Despite its success, RL encounters significant hurdles in complex and partially observable environments, where agents often operate with limited or noisy information (Madani et al., 1999). However, during training, we often have access to extra information that could significantly enhance learning efficiency and performance (Lee et al., 2020a; Chen et al., 2022). For instance, in robotics, while real-world sensor data may be noisy or incomplete, simulation environments typically provide full state observability.

Despite the potential of such privileged information, effectively leveraging it in practice remains a major challenge. One popular strategy to utilize this information is through methods like Imitation Learning (IL) (Hussein et al., 2017), Teacher-Student Learning (TSL), or policy distillation (Czarnecki et al., 2019). In these approaches, a teacher, equipped with privileged information, provides supervision to guide the student’s learning process. However, this strategy introduces its own set of challenges: a teacher with privileged information may impose an unrealistically high-performance standard, making it difficult for the student to effectively imitate. This issue, known as the “impossibly good” teacher (Walsman et al., 2023) or imitation gap (Weihs et al., 2024), can hinder learning and degrade performance. To address this, previous work has sought to integrate environmental rewards into the learning process of the student. One approach is to combine RL with IL (Weihs et al., 2024; Shenfeld et al., 2023a; Nguyen et al., 2023), switching to RL-based training when the teacher becomes inimitable. Another approach modifies environmental rewards based on the teacher through policy distillation (Czarnecki et al., 2019; Walsman et al., 2023). However, such methods diminish the utility of privileged information, often resulting in inefficient use of the teacher’s knowledge.

To better exploit available information, we propose training a “possibly good” teacher. Inspired by Guided Policy Search (GPS) (Levine & Koltun, 2013; Montgomery & Levine, 2016), we introduce Guided Policy Optimization (GPO), a novel framework that trains both the teacher and the student simultaneously while ensuring that the teacher’s policy remains aligned with that of the student. The key insight behind GPO is that by leveraging privileged information during training, the teacher can be trained more effectively while ensuring that its performance is “possibly good,” thus facilitating easier imitation by the student. Theoretically, we show that the student can achieve optimality similar

to direct RL training, mitigating the suboptimality and imitation gaps that often arise from purely teacher-based supervision. We empirically validate our algorithm across various tasks, including didactic examples, challenging continuous control tasks in partially observable, noisy environments within the Brax (Freeman et al., 2021) domain, and in memory-based tasks from the POPGym (Morad et al., 2023) benchmark. GPO shows consistent and significant improvements, underscoring its ability to exploit extra information and deliver robust performance across diverse domains.

2 BACKGROUND

We consider Partially Observable Markov Decision Process (POMDP) (Kaelbling et al., 1998), which is characterized by the tuple $\langle \mathcal{S}, \mathcal{A}, r, \mathcal{P}, \mathcal{O}, \gamma \rangle$. \mathcal{S} represents the set of states, \mathcal{A} the set of actions, r the reward function, \mathcal{P} the transition probability function, \mathcal{O} the partial observation function and γ the discount factor. At each time step t , the agent receives a partial observation $o_t \sim \mathcal{O}(\cdot|s_t)$ for current state $s_t \in \mathcal{S}$. The agent then selects an action $a_t \in \mathcal{A}$ according to o_t or its action-observation history $\tau_t : \{o_0, a_0, o_1, a_1, \dots, o_t\}$. The state transitions to the next state s_{t+1} according to $\mathcal{P}(s_{t+1}|s_t, a_t)$, and the agent receives a reward r_t . The goal for the agent is to find the optimal policy $\pi^* : \tau \rightarrow \Delta(\mathcal{A})$ that maximizes the return, expressed as $\pi^* = \arg \max_{\pi} V_{\pi}$, where $V_{\pi} = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \pi]$ represents cumulative rewards. When full state information s is available during training, we may also define a policy $\mu : s \rightarrow \Delta(\mathcal{A})$ based on privileged information. For clarity, throughout this paper we refer to such privileged training inputs simply as the state s , though in practice they could take other forms. Likewise, we refer to partial observations simply as o , though in practice they may include histories or other derived features.

Finally, we emphasize that in the remainder of this paper, the term “optimal” refers to the student’s optimal policy under partial observability—not the teacher’s optimal policy under privileged information, which is generally unattainable for the student.

2.1 TEACHER-STUDENT LEARNING

Since we consider both training the teacher and student, in this paper, we use the term Teacher-Student Learning (TSL) to broadly refer to Imitation Learning (IL) (Hussein et al., 2017), policy distillation (Czarnecki et al., 2019), and related approaches, as there is no fundamental distinction between them. In TSL, the teacher policy is typically pre-trained using RL or derived from other methods such as a classical controller, which is assumed to effectively accomplish the desired task. The goal is for the teacher to somehow provide supervision to the student in learning to solve the same task.

A straightforward approach to training the agent is to directly supervise the student’s policy π using the teacher’s policy μ , similar to Behavioral Cloning (BC) (Pomerleau, 1991; Torabi et al., 2018):

$$\min_{\pi} \mathbb{E}_{s \sim d_{\mu}} [D_{\text{KL}}(\mu(\cdot|s), \pi(\cdot|s))], \quad (1)$$

where d_{μ} is the distribution of states under the teacher’s policy, and D_{KL} is the Kullback-Leibler (KL) divergence. This objective encourages the student’s policy to mimic the teacher’s policy for the observed states. However, when the teacher’s policy is based on privileged information, the student can only learn the statistical average of the teacher’s actions (Warrington et al., 2020; Weihs et al., 2024), and be strictly suboptimal (Cai et al., 2024). In this paper, we refer to such a teacher as *inimitable*, and we highlight this limitation through two illustrative examples in the next subsection.

2.2 DIDACTIC EXAMPLES

TigerDoor. In the classic TigerDoor problem (Littman et al., 1995), there are two doors with a tiger hidden behind one of them. The possible state s_L (tiger behind the left door) and s_R (tiger behind the right door), with equal probabilities for each, form $\mathcal{S} = \{s_L, s_R\}$. The action set is $\mathcal{A} = \{a_L, a_R, a_I\}$, where a_L and a_R denote opening the left and right doors, respectively, and a_I denotes listening to determine the tiger’s location. The teacher knows the tiger’s location whereas the student can only ascertain it after choosing a_I . The payoff matrix is shown in Table 2. The optimal policy for the teacher is to always choose the correct door without listening, whereas the student’s optimal strategy involves first listening to locate the tiger. Consequently, the student cannot learn the optimal policy through supervision from the teacher, as the teacher never chooses a_I . Under the

teacher’s supervision, the student will only learn to randomly select between a_L and a_R , resulting in an expected reward of 0.5. This scenario poses challenges for the supervised student, as the teacher fails to explore and gather essential information for the learner. We also introduce an alternative version of the problem, TigerDoor-alt (Table 3), which further illustrates the imitation gap, even when no exploratory actions are required. A detailed description is provided in Appendix D.

The current solution to this issue is to incorporate rewards into the student’s learning process. To effectively utilize teacher supervision, there are two kinds of approaches. The first type dynamically adjusts the weight of the supervision loss between the teacher and pure RL training. This allows the algorithm to switch to pure RL when the teacher is deemed inimitable (Weihs et al., 2024; Shenfeld et al., 2023a;b). However, such approaches fail to fully utilize privileged information and may waste the valuable, often expensive, pre-trained teacher. The second kind incorporates teacher supervision into the reward signal, for instance, by using reward shaping via the teacher’s value function (Walsman et al., 2023). However, this supervision is indirect and may require additional learning. Crucially, to the best of our knowledge, none of the existing methods provide theoretical guarantees that teacher supervision will actually be beneficial.

Another research direction attempts to reconstruct privileged information from partial observations. However, such methods require the MDP to be *decodable* (Efroni et al., 2022), which is clearly infeasible in the TigerDoor setting. For a more detailed discussion of related work, see Appendix A.

3 METHOD

We present our Guided Policy Optimization (GPO) framework, which co-trains two entities: the guider and the learner, which we use to differentiate from existing TSL methods. GPO iteratively updates both policies to ensure alignment. We then explore both the theoretical properties and practical implementation of GPO, introducing two variants: GPO-penalty and GPO-clip.

3.1 FROM GPS TO GPO

Unlike direct policy search methods, GPS (Levine & Koltun, 2013; Montgomery & Levine, 2016) does not optimize policy parameters directly. Instead, it introduces an intermediate agent (guider) and employs trajectory optimization to learn a time-varying linear-Gaussian policy, which is then used to train a neural network policy (learner) through supervised learning. Although GPS is a model-based method and is not directly applicable in our setting, its idea of introducing an intermediate agent to guide policy learning can be extended to an RL algorithm in the context of POMDPs.

Specifically, since the guider is only used during training, it can access any type of privileged information. The key requirement is to ensure that the guider is imitable by the learner, which motivates us to introduce the GPO framework, which operates through the following four key steps:

- **Data Collection:** Collect trajectories by executing the guider’s policy, denoted as $\mu^{(k)}$.
- **Guider Training:** Update the guider $\mu^{(k)}$ to $\hat{\mu}^{(k)}$ according to RL objective $V_{\mu^{(k)}}$.
- **Learner Training:** Update the learner to $\pi^{(k+1)}$ by minimizing the distance $D(\pi, \hat{\mu}^{(k)})$.
- **Guider Backtracking:** Set $\mu^{(k+1)}(\cdot|s) = \pi^{(k+1)}(\cdot|o)$ for all states s before the next iteration.

In the learner training step, $D(\pi, \mu)$ can be any Bregman divergence. For this work, we use the KL divergence, weighted by the state distribution d_μ . GPO iterates these steps until convergence, applying standard RL to train the guider while the learner seeks to mimic the guider’s behavior. If the learner struggles due to discrepancies in observation spaces, the backtracking step adjusts the guider’s policy to mitigate the imitation gap.

The comparison between TSL and GPO is illustrated in Fig. 1. Several key differences between the two frameworks exist. First, the teacher in TSL is typically provided or trained independently from the student, while in GPO, the guider and learner are trained together. Second, TSL typically allows the student to interact with the environment, whereas GPO only uses a guider, enabling more effective trajectory collection due to the behavioral policy being conditioned on privileged information. Lastly, and most importantly, TSL does not use the student to constrain the teacher. This means that if the teacher is too advanced for the student, the student will struggle to learn from the teacher.

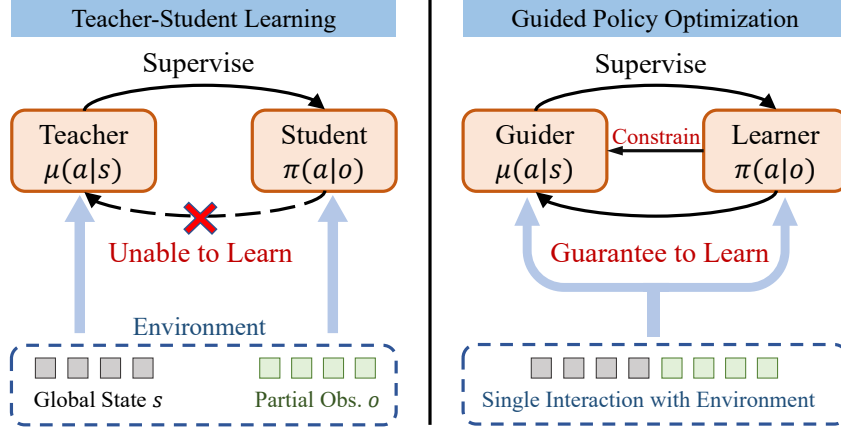


Figure 1: The comparison between TSL and GPO.

In contrast, GPO utilizes backtracking to guarantee the learner can effectively learn from the guider. This is demonstrated by the following proposition:

Proposition 1. *If the guider’s policy is updated using policy mirror descent in each GPO iteration:*

$$\hat{\mu} = \arg \min \{ -\eta_k \langle \nabla V(\mu^{(k)}), \mu \rangle + D_{\mu^{(k)}}(\mu, \mu^{(k)}) \},$$

where η_k is the step size. Then the learner’s policy update follows a constrained policy mirror descent:

$$\pi^{(k+1)} = \arg \min_{\pi \in \Pi} \{ -\eta_k \langle \nabla V(\pi^{(k)}), \pi \rangle + D_{\pi^{(k)}}(\pi, \pi^{(k)}) \}.$$

Proof. See Appendix B. □

Here, we assume that the guider μ has access to an unlimited policy class, while the learner π is constrained to a limited policy class Π for simplicity. Policy mirror descent (Tomar et al., 2020; Xiao, 2022) is a general family of algorithms that encompasses a wide range of fundamental methods in RL, including trust-region algorithms like TRPO (Schulman et al., 2015a) and PPO (Schulman et al., 2017). This proposition shows that, despite the learner not directly interacting with the environment, the GPO update for the learner can be viewed as a standard RL update. Specifically, if we use trust-region RL algorithms for the guider, the update for the learner’s policy inherits the key properties, such as policy improvement (Schulman et al., 2015a). This suggests that GPO can effectively address challenges in TSL, such as dealing with a suboptimal teacher or the imitation gap, while still framing the learner’s policy as being supervised by the guider. In Appendix D, we provide an intuitive example illustrating how GPO can achieve optimal in the TigerDoor-alt problem.

Given that GPO mirrors direct RL for the learner, one may ask: **What are GPO’s key advantages?** The main benefit lies in leveraging additional information while simplifying learning. Since policy gradients suffer from high variance especially under partial observability, GPO splits learning into two phases: the guider with privileged information handles complex RL gradients, while the partial observable learner is trained via an easier supervised learning, reducing variance and complexity. For instance, to train robustness to noisy observations, GPO can train the guider on clean inputs and supervise the learner with noisy ones, resulting in a more stable and effective learning process.

3.2 GPO-PENALTY

This section introduces a straightforward implementation of the GPO framework using KL-divergence as a penalty for the guider, which we refer to as GPO-penalty. Specifically, in step 2 of GPO, we use PPO as the underlying trust-region algorithm. The corresponding loss for the guider’s policy is as follows¹:

$$\mathcal{L}_1(\mu) = -\mathbb{E} \left[\min \left(\rho^\mu A^\beta(s, a), \rho_{clip}^\mu A^\beta(s, a) \right) \right], \quad (2)$$

¹We omit subscripts for expectations in the remainder of the paper, as all samples are drawn from the distribution induced by the behavioral policy $\beta = \mu_{old}$.

where $\rho^\mu = \mu(a|s)/\beta(a|s)$, $\rho_{clip}^\mu = \text{clip}(\rho^\mu, 1 - \epsilon, 1 + \epsilon)$ and β denotes the behavioral policy. The advantage $A^\beta(s, a)$ is estimated using the Generalized Advantage Estimation (GAE) (Schulman et al., 2015b) with the value function $V(s)$ trained via discounted reward-to-go.

In step 3, since finding the exact minimizer of the distance measure is computationally prohibitive, we use gradient descent to minimize the BC objective:

$$\mathcal{L}_2(\pi) = \mathbb{E}[\text{D}_{\text{KL}}(\mu(\cdot|s), \pi(\cdot|o))]. \quad (3)$$

Similarly, in step 4, we backtrack the guider’s policy using the same BC loss:

$$\mathcal{L}_3(\mu) = \mathbb{E}[\text{D}_{\text{KL}}(\mu(\cdot|s), \pi(\cdot|o))]. \quad (4)$$

A key insight in GPO is that exact backtracking of the guider’s policy is unnecessary—it’s sufficient to keep the guider within a imitable region relative to the learner. The learner may fail to follow the guider either because the guider is inimitable or just because the guider learns faster, the latter being common due to inexact gradient updates. In such cases, aggressive backtracking can be harmful. Keeping the guider slightly ahead also allows it to collect better trajectories, as discussed in Section 4.4. To maintain this balance, we introduce a coefficient α that modulates the guider’s loss as

$$\mathcal{L}(\mu) = \mathcal{L}_1(\mu) + \alpha \mathcal{L}_3(\mu), \quad (5)$$

where α is adapted based on the distance $\mathcal{L}_3(\mu)$ relative to a threshold d , using a constant scaling factor k :

$$\alpha = k\alpha \text{ if } \mathcal{L}_3(\mu) > kd, \quad \alpha/k \text{ if } \mathcal{L}_3(\mu) < d/k. \quad (6)$$

This scheme is analogous to the KL-penalty adjustment in PPO-penalty (Schulman et al., 2017), where the penalty coefficient adjusts based on the relationship between the KL divergence and a predefined threshold.

Another key aspect is compensating for the learner’s policy improvement, as we replace strict backtracking with a KL constraint. While it is possible to set a very small d_{targ} , this would inefficiently inflate α , hindering the guider’s training. Notably, Proposition 1 implies that applying GPO with PPO is effectively equivalent to applying PPO directly to the learner. Consequently, we can concurrently train the learner’s policy using PPO during the GPO iterations. As a result, we introduce an additional objective for the learner’s policy:

$$\mathcal{L}_4(\pi) = -\mathbb{E}\left[\min\left(\rho^\pi A^\beta(s, a), \rho_{clip}^\pi A^\beta(s, a)\right)\right], \quad (7)$$

where $\rho^\pi = \pi(a|o)/\beta(a|s)$. Considering that the behavioral policy is from the guider, to validate this update, we introduce the following proposition:

Proposition 2. For policy π , μ , β and all states s , suppose $D_{TV}(\mu(\cdot|s), \beta(\cdot|s)) \lesssim \epsilon/2$, then we have

$$\mathbb{E}_{a \sim \beta}[|1 - \rho^\pi(s, a)|] \lesssim \epsilon + \sqrt{2d_{\text{targ}}}.$$

Proof. See Appendix B. □

The assumption on total variation distance is justified by the PPO update of the guider’s policy (Appendix B). This proposition implies that when d_{targ} is small, the behavioral policy closely matches the learner’s policy, allowing valid sample reuse for learner training.

Finally, we define the merged learner objective for the learner as:

$$\mathcal{L}(\pi) = \alpha \mathcal{L}_4(\pi) + \mathcal{L}_2(\pi), \quad (8)$$

where the coefficient α from equation 6 is applied to the RL term. This mechanism compensates when the learner struggles to follow the guider. If the learner is able to fully track the guider, α approaches zero, allowing the guider to directly lead the learner to the optimal policy without requiring an additional RL objective. When the learner cannot keep pace, the RL objective aids in the learner’s training.

3.3 GPO-CLIP

In this section, we introduce a slightly modified implementation of the GPO framework, which we refer to as GPO-clip. The key principle is that an effective guider should remain at the boundary of the learner’s imitable region: if the guider is too far ahead, the learner struggles to follow; if too close, the guider’s ability to provide effective supervision and better trajectory diminishes. To achieve this balance, the guider should halt updates when it moves too far ahead and avoid backtracking when it is already sufficiently close.

We propose two key modifications to the GPO-penalty algorithm introduced in the previous subsection. First, inspired by PPO-clip, we replace the clip function ρ_{clip}^μ in equation 2 with the following double-clip function:

$$\rho_{clip}^{\mu,\pi} = \text{clip}\left(\text{clip}\left(\frac{\mu(a|s)}{\pi(a|o)}, 1 - \delta, 1 + \delta\right) \cdot \frac{\pi(a|o)}{\beta(a|s)}, 1 - \epsilon, 1 + \epsilon\right). \quad (9)$$

This formulation introduces an additional inner clipping step, which halts the guider’s updates under two conditions: (1) $A^\beta(s, a) > 0$ and $\mu(a|s) > \pi(a|o)(1 + \delta)$, (2) $A^\beta(s, a) < 0$ and $\mu(a|s) < \pi(a|o)(1 - \delta)$. Considering that the positive (negative) advantage indicates that $\mu(a|s)$ is set to increase (decrease), the double-clip function prevents further movement away from π when μ is already distant.

It is important to note that, unlike PPO where PPO-clip can completely replace the KL-penalty term, this is not the case in GPO. In PPO, the ratio $\rho^\pi(s, a)$ starts at 1 at the beginning of each epoch, ensuring that the clipped ratio keeps π near the behavioral policy. In GPO, however, the gap between $\pi(a|s)$ and $\mu(a|o)$ may accumulate over multiple updates if the learner fails to keep up with the guider. The double-clip function equation 9 alone is insufficient to bring $\pi(a|o)$ back into the δ region once it has strayed too far. To address this, we introduce a mask on the backtracking loss, defined as:

$$m(s, a) = \mathbb{I}\left(\frac{\mu(a|s)}{\pi(a|o)} \notin (1 - \delta, 1 + \delta)\right), \quad (10)$$

where \mathbb{I} is the indicator function. This mask replaces the adaptive coefficient α of GPO-penalty, selectively applying the backtracking penalty only when $\mu(a|o)$ drifts outside the δ region. Policies that remain close to each other are left unaffected, preventing unnecessary backtracking.

Additionally, given that both the guider and learner are solving the same task, their policies should exhibit structural similarities. To leverage this, we allow the guider and learner to share a single policy network. To distinguish between guider and learner inputs, we define a unified input format: the input to the guider’s policy is defined as $o_g = [s, o, 1]$, where s is the state, o is the partial observation, and the scalar 1 serves as an indicator; the learner’s input is defined as $o_l = [\vec{0}, o, 0]$, where $\vec{0}$ is a zero vector with the same dimensionality as s , indicating that the learner has access only to the partial observation o . This approach is applied to both GPO-penalty and GPO-clip, and the update for the shared policy network with parameters θ is as follows:

$$L_{\text{GPO-penalty}}(\theta) = \mathbb{E}\left[-\min\left(\rho^{\mu_\theta} A^\beta(o_g, a), \rho_{clip}^{\mu_\theta} A^\beta(o_g, a)\right) + \alpha \text{D}_{\text{KL}}(\mu_\theta(\cdot|o_g) || \pi_{\hat{\theta}}(\cdot|o_l))\right. \\ \left.- \alpha \min\left(\rho^{\pi_\theta} A^\beta(o_l, a), \rho_{clip}^{\pi_\theta} A^\beta(o_l, a)\right) + \text{D}_{\text{KL}}(\mu_{\hat{\theta}}(\cdot|o_g) || \pi_\theta(\cdot|o_l))\right], \quad (11)$$

$$L_{\text{GPO-clip}}(\theta) = \mathbb{E}\left[-\min\left(\rho^{\mu_\theta} A^\beta(o_g, a), \rho_{clip}^{\mu_\theta, \pi_{\hat{\theta}}} A^\beta(o_g, a)\right) + m(s, a) \text{D}_{\text{KL}}(\mu_\theta(\cdot|o_g) || \pi_{\hat{\theta}}(\cdot|o_l))\right. \\ \left.- \alpha \min\left(\rho^{\pi_\theta} A^\beta(o_g, a), \rho_{clip}^{\pi_\theta} A^\beta(o_g, a)\right) + \text{D}_{\text{KL}}(\mu_{\hat{\theta}}(\cdot|o_g) || \pi_\theta(\cdot|o_l))\right], \quad (12)$$

where $\hat{\theta}$ denotes a stop-gradient operation on the parameters, and α for GPO-clip is a fixed parameter. The detailed algorithms are summarized in Appendix C.

4 EXPERIMENTS

In this section, we evaluate the empirical performance of GPO across various domains. For baselines, we consider two types of approaches for utilizing teacher supervision. The first type involves training both the teacher and student simultaneously. A summary of their main characteristics is provided in Table 1. Among these, **GPO-naive** refers to GPO-penalty without the RL auxiliary loss. **PPO-asym** directly trains the learner using PPO,

with the learner’s value function receiving o_g as input. **PPO+BC** trains the teacher with PPO, while the learner is trained via direct BC from the teacher. **ADVISOR-co** is a modification of ADVISOR (Weihs et al., 2024), and **A2D** is based on the work by (Warrington et al., 2020). The second type involves training the teacher first, followed by the application of TSL methods. These include **Dagger** (Ross et al., 2011), **PPO+BC-t**, **ADVISOR**, **ELF** (Walsman et al., 2023), and **ELF-*asym***, where ELF is a policy distillation method that utilizes reward shaping to provide supervision to the student, and ELF-*asym* is a variant that uses an asymmetric value function. Further details about these algorithms can be found in Appendix E.1.

Table 1: Co-training algorithms.

Algorithm	Train μ	Behavioral policy	Train π	Value function
PPO	-	$\pi(a o_l)$	PPO	$V(o_l)$
PPO- <i>asym</i>	-	$\pi(a o_l)$	PPO	$V(o_g)$
PPO+BC	PPO	$\mu(a o_g)$	BC	$V(o_g)$
A2D	PPO	$\pi(a o_l)$	BC	$V(o_l)$
ADVISOR-co	PPO	$\pi(a o_l)$	BC+PPO	$V(o_l)$
GPO-naive	PPO	$\mu(a o_g)$	BC	$V(o_g)$
GPO-penalty	PPO	$\mu(a o_g)$	BC+PPO	$V(o_g)$
GPO-clip	PPO	$\mu(a o_g)$	BC+PPO	$V(o_g)$
GPO-ablation	PPO	$\mu(a o_g)$	PPO	$V(o_g)$

4.1 DIDACTIC TASKS

Table 2: TigerDoor

$s \backslash a$	a_L	a_R	a_l
s_L	1	0	-0.1
s_R	0	1	-0.1

Table 3: TigerDoor-alt

$s \backslash a$	a_L	a_R
s_L	2	0
s_R	0	1

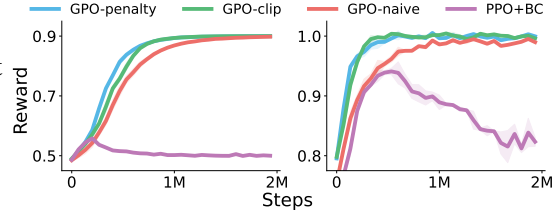


Figure 2: Results on TigerDoor (left) and TigerDoor-alt (right).

We begin by evaluating our algorithm on two didactic problems introduced in Section 2.2. As shown in Fig. 2, direct cloning of the guider’s policy converges to a suboptimal solution, as expected. In contrast, all variants of GPO achieve optimal performance on these tasks. Although applying RL directly to the learner easily leads to optimal solutions, it is important to note that GPO-naive achieves optimality purely through supervised learning. This result verifies the optimality guarantee of the GPO framework described in Proposition 1, suggesting that a guider constrained within the learner’s imitable region can provide effective supervision, even with asymmetric information. Moreover, comparing GPO-naive to GPO-penalty and GPO-clip reveals that the introduction of direct RL training for the learner accelerates learning.

4.2 CONTINUOUS CONTROL TASKS IN BRAX

In this subsection, we present the results of our algorithms and baselines on several continuous control tasks in the Brax domain. To transform these tasks into a POMDP setting, we remove the velocity information of all joints, and add varying levels of noise to the observations. The guider has access to full, noiseless information, while the learner operates with partial and noisy inputs. For more details, please refer to Appendix E.

The results are shown in Fig. 3, where the performance hierarchy is generally: GPO-clip > GPO-penalty > PPO-*asym* > GPO-naive > other baselines. It is important to note that, even without factoring in the cost of training the teacher (which is nearly the same as training the GPO algorithm from scratch), methods that rely on a pre-trained privileged teacher perform well only in the *Halfcheetah* and *Swimmer* tasks. Furthermore, the performance of these methods declines rapidly as the noise scale increases. This occurs because, when the pre-trained teacher becomes too skilled for the student, it provides little to no useful supervision, and may even have a negative impact on learning.

For co-training approaches, we have the following key observations: First, the superior performance of GPO-clip and GPO-penalty compared to the base algorithm PPO shows that this framework can effectively utilize additional information during training to facilitate the learner training. Second, comparing GPO-naive to GPO-penalty and GPO-clip, we see that introducing RL training for the learner improves performance. Third, the comparison between PPO+BC and GPO-naive highlights

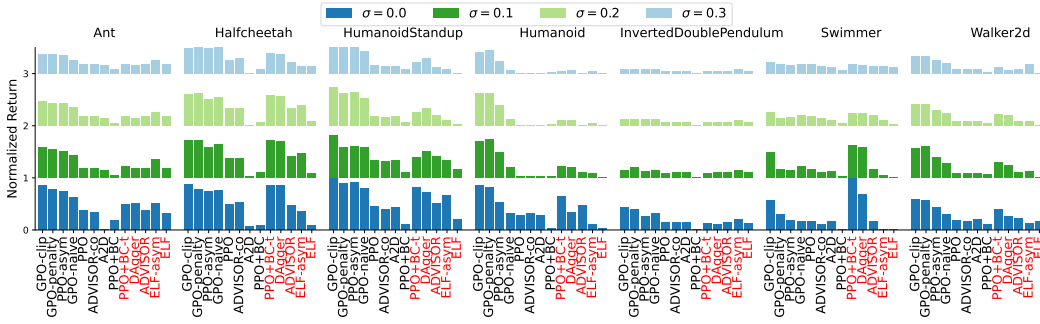


Figure 3: Comparison of GPO and baselines on the Brax domain, where σ represents the scale of Gaussian noise added to the observations. The performance on each task is normalized to $[0, 1]$ using the performance of the corresponding pre-trained teacher as a reference. Algorithms highlighted in red are supervised by their corresponding pre-trained teacher.

the necessity of backtracking. If the guider is not constrained to the learner, the guider’s supervision may negatively influence the performance. Last, other baselines such as ADVISOR failed to utilize the privileged teacher such that it degenerates into pure PPO.

In summary, our method consistently outperforms the baselines, demonstrating its effectiveness in solving noisy and partially observable continuous control tasks. Additional experiments including L2T-RL (Wu et al., 2024), TGRL (Shenfeld et al., 2023b) and RMA (Kumar et al., 2021) are provided in Appendix E.4.

4.3 MEMORY-BASED TASKS IN POPGYM

Since using memory models to deal with POMDP is a common practice, we evaluate GPO in POPGym to show whether the algorithm can effectively address memory-based tasks. The tasks include card and board games where agents must recall previous observations to extract useful information for decision-making. For these tasks, the guider’s observation is designed to include the critical information needed to remember, theoretically minimizing the imitation gap as long as the memory model can store the necessary information. Although in practice, memory models struggle to retain all information, especially in complex tasks, this setup allows us to use a larger KL threshold or clipping parameter, enabling the guider to explore further and provide more valuable supervision. Further details on the experimental settings are provided in Appendix E.

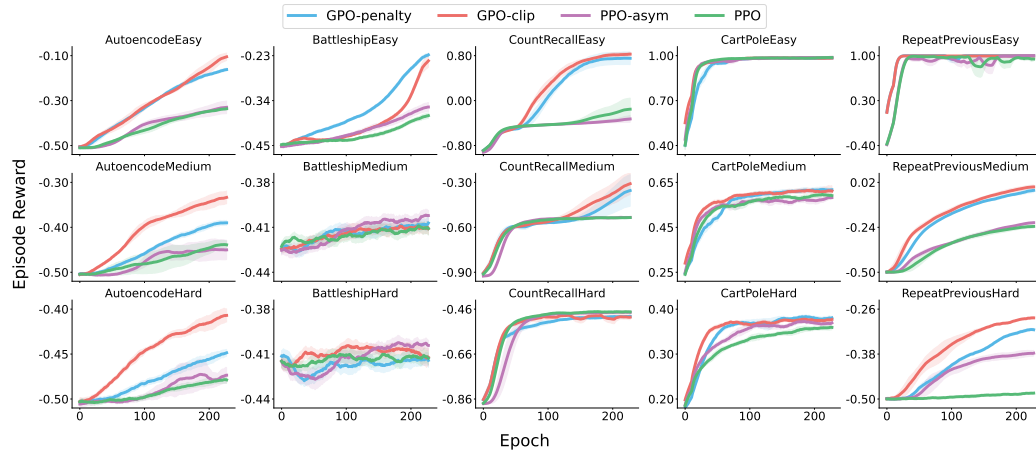


Figure 4: The results of GPO-clip, GPO-penalty, PPO-asym, and PPO on 15 POPGym tasks.

Fig. 4 shows the results on 15 POPGym tasks, where we compare GPO-penalty and GPO-clip to PPO-asym and PPO. The general conclusion mirrors the results from the previous subsection, where GPO-clip typically outperforms GPO-penalty, followed by PPO-asym and PPO. Key insights include: First, the superior performance of GPO-penalty indicates that the ability of the guider to explore

further without diverging too much from the learner proves valuable in these memory-based tasks. Second, while PPO-asym outperforms PPO, its performance improvement is less pronounced here than in the Brax domain, suggesting that asymmetric value function may not be very helpful for memory tasks. Third, although neither GPO-penalty nor GPO-clip exhibits superior performance in tasks like *BattleshipMedium* and *CountRecallHard*, this is due to we use the same parameter across all tasks, and performance could be improved as we show in the next section.

Overall, our methods demonstrate strong performance across the majority of tasks, providing an effective solution for memory-based problems.

4.4 ABLATIONS AND DISCUSSIONS

In this section, we dive deeper into GPO’s performance through ablations and further discussions.

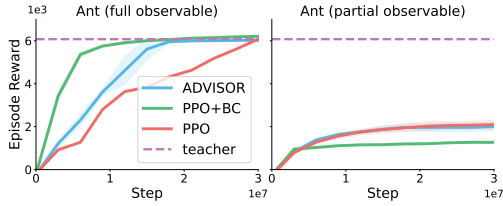


Figure 5: ADVISOR and PPO+BC with a pre-trained teacher.

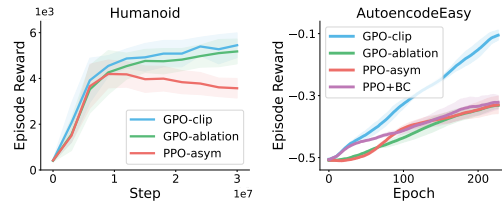
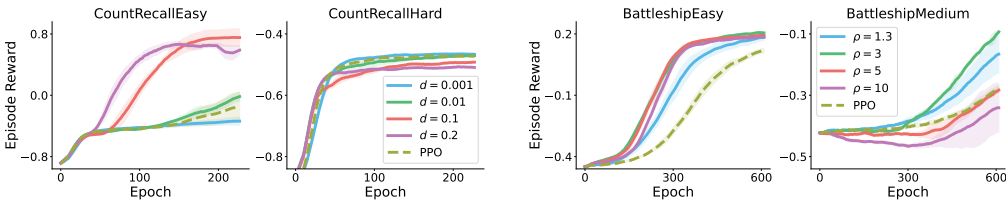


Figure 6: Ablation studies.

Why does training a teacher first and applying TSL methods often fail? A representative example is the TigerDoor problem, where a pre-trained teacher provides minimal to no effective supervision for the student. Recent TSL approaches, such as ADVISOR and TGRL, address the challenge of an overly optimal teacher by reverting to pure RL, thereby bypassing uninformative or misleading supervision. As shown in Fig. 5, although ADVISOR and PPO+BC perform well in the fully observable *Ant* task where the teacher is trained, it degenerates into PPO in the partially observable *Ant* task since the teacher is found inimitable.

Why do GPO outperform other baselines? We attribute the superior performance to two factors: effective RL training of the learner, and effective supervision from the guider. The benefit of RL training is shown in Fig. 6(left), where GPO-ablation (GPO-penalty without supervision, as described in Table 1) outperforms PPO-asym on the *Humanoid* task. Although both use similar objectives, GPO-ablation uses data collected by the guider, indicating that a better behavior policy improves learning efficiency. The effectiveness of the supervision comes from the guider being constrained to the imitable region while still learning rapidly. In Fig. 6(right), with the learner trained purely by supervision (GPO-clip with RL disabled), GPO-clip outperforms GPO-ablation, PPO+BC, and PPO-asym. This shows that in memory-intensive tasks, supervision is more beneficial than RL. Since PPO+BC performs poorly in noisy tasks in Section 4.2 but comparably to PPO-asym here, we can also infer that supervision plays a particularly important role in these tasks. Moreover, GPO-clip’s strong performance over PPO+BC—despite both using pure supervision—highlights the importance of constraining the guider to a policy the learner can follow.



(a) GPO-penalty with different KL-threshold.

(b) GPO-clip with different clip parameters.

Figure 7: The results of GPO-penalty and GPO-clip with different hyperparameters. The clip parameter ρ is defined in Appendix E.2.

When does GPO fail? GPO can fail when the guider learns too slowly, often due to inadequate information. Another failure mode arises from poorly tuned KL thresholds (clip parameters). For instance,

in the *CountRecallHard* task from POPGym, both GPO variants underperform compared to PPO and PPO-asym. As shown in Fig. 7, larger KL thresholds help in simple tasks like *CountRecallEasy* and *BattleshipEasy*, but hurt performance in harder ones like *CountRecallHard* and *BattleshipMedium*. This is because challenging tasks strain memory models like GRU—when GRU fails to retain key information, the learner cannot follow the guider. In such cases, a large KL threshold pushes the guider beyond the learner’s reachable region, causing an unrecoverable imitation gap.

5 CONCLUSION AND FUTURE WORK

In this paper, we introduce GPO, a method designed to leverage additional information in POMDPs during training. Our experimental results demonstrate that the proposed algorithm effectively addresses noisy and memory-based partially observable tasks, offering a novel approach to utilizing auxiliary information for more efficient learning. Future work could explore extending guided policy optimization to the multi-agent setting, where agents often have access to global information during training but are constrained to local observations during execution.

REFERENCES

- OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Józefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider, Szymon Sidor, Josh Tobin, Peter Welinder, Lilian Weng, and Wojciech Zaremba. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020. doi: 10.1177/0278364919887447.
- Andrea Baisero and Christopher Amato. Unbiased asymmetric reinforcement learning under partial observability. *arXiv preprint arXiv:2105.11674*, 2021.
- Mayank Bansal, Alex Krizhevsky, and Abhijit Ogale. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. *arXiv preprint arXiv:1812.03079*, 2018.
- Mohak Bhardwaj, Sanjiban Choudhury, and Sebastian Scherer. Learning heuristic search via imitation. In *Conference on Robot Learning*, pp. 271–280. PMLR, 2017.
- Yang Cai, Xiangyu Liu, Argyris Oikonomou, and Kaiqing Zhang. Provable partially observable reinforcement learning with privileged information. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 63790–63857. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/74d188c51d97fcfb0269f584d6a53b7-Paper-Conference.pdf.
- Kai-Wei Chang, Akshay Krishnamurthy, Hal Daumé, III, and John Langford. Learning to search better than your teacher. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2058–2066, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/changb15.html>.
- Tao Chen, Jie Xu, and Pulkit Agrawal. A system for general in-hand object re-orientation. In Aleksandra Faust, David Hsu, and Gerhard Neumann (eds.), *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pp. 297–307. PMLR, 08–11 Nov 2022. URL <https://proceedings.mlr.press/v164/chen22a.html>.
- Sanjiban Choudhury, Ashish Kapoor, Gireeja Ranade, Sebastian Scherer, and Debadeepta Dey. Adaptive information gathering via imitation learning. *arXiv preprint arXiv:1705.07834*, 2017.
- Wojciech M. Czarnecki, Razvan Pascanu, Simon Osindero, Siddhant Jayakumar, Grzegorz Swirszcz, and Max Jaderberg. Distilling policy distillation. In Kamalika Chaudhuri and Masashi Sugiyama (eds.), *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 1331–1340. PMLR, 16–18 Apr 2019. URL <https://proceedings.mlr.press/v89/czarnecki19a.html>.

- Pim De Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning. *Advances in neural information processing systems*, 32, 2019.
- Yonathan Efroni, Chi Jin, Akshay Krishnamurthy, and Sobhan Miryoosefi. Provable reinforcement learning with a short-term memory. In *International Conference on Machine Learning*, pp. 5832–5850. PMLR, 2022.
- Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax - a differentiable physics engine for large scale rigid body simulation, 2021. URL <http://github.com/google/brax>.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- Maximilian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, and Shimon Whiteson. Deep variational reinforcement learning for pomdps. In *International conference on machine learning*, pp. 2117–2126. PMLR, 2018.
- Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134, 1998. ISSN 0004-3702. doi: [https://doi.org/10.1016/S0004-3702\(98\)00023-X](https://doi.org/10.1016/S0004-3702(98)00023-X). URL <https://www.sciencedirect.com/science/article/pii/S000437029800023X>.
- Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. Rma: Rapid motor adaptation for legged robots. *arXiv preprint arXiv:2107.04034*, 2021.
- John Lambert, Ozan Sener, and Silvio Savarese. Deep learning under privileged information using heteroscedastic dropout. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8886–8895, 2018.
- Jonathan Lee, Alekh Agarwal, Christoph Dann, and Tong Zhang. Learning in pomdps is sample-efficient with hindsight observability. In *International Conference on Machine Learning*, pp. 18733–18773. PMLR, 2023.
- Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning quadrupedal locomotion over challenging terrain. *Science Robotics*, 5(47):eabc5986, 2020a. doi: 10.1126/scirobotics.abc5986. URL <https://www.science.org/doi/abs/10.1126/scirobotics.abc5986>.
- Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning quadrupedal locomotion over challenging terrain. *Science robotics*, 5(47):eabc5986, 2020b.
- Sergey Levine and Vladlen Koltun. Guided policy search. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 1–9, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/levine13.html>.
- Michael L. Littman, Anthony R. Cassandra, and Leslie Pack Kaelbling. Learning policies for partially observable environments: scaling up. In *Proceedings of the Twelfth International Conference on Machine Learning*, ICML’95, pp. 362–370, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603778.
- Qinghua Liu, Alan Chung, Csaba Szepesvári, and Chi Jin. When is partially observable reinforcement learning not scary? In *Conference on Learning Theory*, pp. 5175–5220. PMLR, 2022.

- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob Foerster, Satinder Singh, and Feryal Behbahani. Structured state space models for in-context reinforcement learning. *arXiv preprint arXiv:2303.03982*, 2023.
- Omid Madani, Steve Hanks, and Anne Condon. On the undecidability of probabilistic planning and infinite-horizon partially observable markov decision problems. *Aaai/iaai*, 10(315149.315395), 1999.
- Lingheng Meng, Rob Gorbet, and Dana Kulić. Memory-based deep reinforcement learning for pomdps. In *2021 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pp. 5619–5626. IEEE, 2021.
- William H Montgomery and Sergey Levine. Guided policy search via approximate mirror descent. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/a00e5eb0973d24649a4a920fc53d9564-Paper.pdf.
- Steven Morad, Ryan Kortvelesy, Matteo Bettini, Stephan Liwicki, and Amanda Prorok. POPGym: Benchmarking partially observable reinforcement learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=chDrutUTs0K>.
- Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Icml*, volume 99, pp. 278–287, 1999.
- Hai Nguyen, Andrea Baisero, Dian Wang, Christopher Amato, and Robert Platt. Leveraging fully observable policies for learning under partial observability. *arXiv preprint arXiv:2211.01991*, 2022.
- Hai Huu Nguyen, Andrea Baisero, Dian Wang, Christopher Amato, and Robert Platt. Leveraging fully observable policies for learning under partial observability. In Karen Liu, Dana Kulic, and Jeff Ichnowski (eds.), *Proceedings of The 6th Conference on Robot Learning*, volume 205 of *Proceedings of Machine Learning Research*, pp. 1673–1683. PMLR, 14–18 Dec 2023. URL <https://proceedings.mlr.press/v205/nguyen23a.html>.
- Lerrel Pinto, Marcin Andrychowicz, Peter Welinder, Wojciech Zaremba, and Pieter Abbeel. Asymmetric actor critic for image-based robot learning. *RSS*, 2018.
- Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Haozhi Qi, Ashish Kumar, Roberto Calandra, Yi Ma, and Jitendra Malik. In-hand object rotation via rapid motor adaptation. In *Conference on Robot Learning*, pp. 1722–1732. PMLR, 2023.
- Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 627–635. JMLR Workshop and Conference Proceedings, 2011.
- Sasha Salter, Dushyant Rao, Markus Wulfmeier, Raia Hadsell, and Ingmar Posner. Attention-privileged reinforcement learning. In *Conference on Robot Learning*, pp. 394–408. PMLR, 2021.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. *CoRR*, abs/1502.05477, 2015a. URL <http://arxiv.org/abs/1502.05477>.

- John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *CoRR*, abs/1506.02438, 2015b. URL <https://api.semanticscholar.org/CorpusID:3075448>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- Younggyo Seo, Junsu Kim, Stephen James, Kimin Lee, Jinwoo Shin, and Pieter Abbeel. Multi-view masked world models for visual robotic manipulation. In *International Conference on Machine Learning*, pp. 30613–30632. PMLR, 2023.
- Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 1134–1141. IEEE, 2018.
- Idan Shenfeld, Zhang-Wei Hong, Aviv Tamar, and Pulkit Agrawal. TGRL: An algorithm for teacher guided reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 31077–31093. PMLR, 23–29 Jul 2023a. URL <https://proceedings.mlr.press/v202/shenfeld23a.html>.
- Idan Shenfeld, Zhang-Wei Hong, Aviv Tamar, and Pulkit Agrawal. Tgrl: An algorithm for teacher guided reinforcement learning. In *International Conference on Machine Learning*, pp. 31077–31093. PMLR, 2023b.
- Jialin Song, Ravi Lanka, Albert Zhao, Yisong Yue, and Masahiro Ono. Learning to search via self-imitation. *CoRR*, abs/1804.00846, 2018. URL <http://arxiv.org/abs/1804.00846>.
- Jialin Song, Ravi Lanka, Yisong Yue, and Masahiro Ono. Co-training for policy learning. In Ryan P. Adams and Vibhav Gogate (eds.), *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine Learning Research*, pp. 1191–1201. PMLR, 22–25 Jul 2020. URL <https://proceedings.mlr.press/v115/song20b.html>.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.
- Voot Tangkaratt, Nontawat Charoenphakdee, and Masashi Sugiyama. Robust imitation learning from noisy demonstrations. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 298–306. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/tangkaratt21a.html>.
- Manan Tomar, Lior Shani, Yonathan Efroni, and Mohammad Ghavamzadeh. Mirror descent policy optimization. *CoRR*, abs/2005.09814, 2020. URL <https://arxiv.org/abs/2005.09814>.
- Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*, 2018.
- Vladimir Vapnik and Akshay Vashist. A new learning paradigm: Learning using privileged information. *Neural networks*, 22(5-6):544–557, 2009.
- Aaron Walsman, Muru Zhang, Sanjiban Choudhury, Ali Farhadi, and Dieter Fox. Impossibly good experts and how to follow them. In *International Conference on Learning Representations*, 2023. URL <https://api.semanticscholar.org/CorpusID:259267611>.
- Andrew Warrington, Jonathan Wilder Lavington, A. Scibior, Mark W. Schmidt, and Frank D. Wood. Robust asymmetric learning in pomdps. In *International Conference on Machine Learning*, 2020. URL <https://api.semanticscholar.org/CorpusID:229923742>.

Luca Weihs, Unnat Jain, Iou-Jen Liu, Jordi Salvador, Svetlana Lazebnik, Aniruddha Kembhavi, and Alexander Schwing. Bridging the imitation gap by adaptive insubordination. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS '21*, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713845393.

Feiyang Wu, Zhaoyuan Gu, Ye Zhao, and Anqi Wu. Learn to teach: Improve sample efficiency in teacher-student learning for sim-to-real transfer, 2024. URL <https://arxiv.org/abs/2402.06783>.

Lin Xiao. On the convergence rates of policy gradient methods. *Journal of Machine Learning Research*, 23(282):1–36, 2022. URL <http://jmlr.org/papers/v23/22-0056.html>.

Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lusong Li, Li Shen, Xiaodong He, Jing Jiang, and Yuhui Shi. Embodied multi-modal agent trained by an llm from a parallel textworld, 2024. URL <https://arxiv.org/abs/2311.16714>.

Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of PPO in cooperative multi-agent games. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.

THE USE OF LARGE LANGUAGE MODELS (LLMs)

LLMs are used to polish the paper writing.

A RELATED WORKS

Although leveraging historical information has proven effective for solving POMDPs (Igl et al., 2018; Meng et al., 2021; Liu et al., 2022), additional information—often available in simulators during training—can be exploited to further aid learning. Leveraging additional information to accelerate learning in POMDPs has been explored across various frameworks and application domains (Vapnik & Vashist, 2009; Lambert et al., 2018; Lee et al., 2023). A prominent line of research focuses on Imitation Learning (IL), where expert knowledge, often equipped with extra information, significantly enhances performance in practical domains like autonomous driving (Bansal et al., 2018; De Haan et al., 2019) and robot navigation and planning (Choudhury et al., 2017; Bhardwaj et al., 2017). However, traditional IL methods such as Behavioral Cloning (BC) (Pomerleau, 1991; Torabi et al., 2018) and DAgger (Ross et al., 2011) often lead to sub-optimal solutions in scenarios requiring active information gathering by the agent (Pinto et al., 2018; Warrington et al., 2020).

To overcome these limitations, recent research has focused on hybrid approaches that integrate RL with IL, often in the context of policy distillation (Czarnecki et al., 2019). For instance, (Nguyen et al., 2022) modifies Soft Actor Critic (SAC) (Haarnoja et al., 2018) by replacing the entropy term with a divergence measure between agent and expert policies at each visited state. Similarly, (Weihs et al., 2024) introduces a balancing mechanism between BC and RL training, adjusting based on the agent’s ability to mimic the expert. Additionally, (Walsman et al., 2023) applies potential-based reward shaping (Ng et al., 1999) using the expert’s value function to guide the agent’s policy gradient, while (Shenfeld et al., 2023b) augments entropy in SAC to blend task reward with expert guidance, where the balance is based on the agent’s performance relative to a reward-only learner.

Despite these advances, expert-driven approaches often assume access to a reliable expert, which may not be feasible when only supplementary information is available. This has led to a growing body of work on co-training approaches where the expert and agent are learned jointly, with the expert conditioned on additional information. For example, (Salter et al., 2021) proposes training separate policies for the agent and expert using spatial attention for image-based RL, aligning attention mechanisms through shared experiences. (Song et al., 2020) co-trains two policies, each conditioned on different information, and selects the most successful rollouts from both policies to guide subsequent learning via RL or IL. (Warrington et al., 2020) further develops this idea in adaptive asymmetric DAgger (A2D), where the expert is continuously refined through RL while supervising the agent. (Wu et al., 2024) also co-trains a teacher and a student, using the experience collected by teacher to apply RL and BC for the student. Beyond expert-based methods, a complementary

approach involves embedding supplementary information directly into the value function within the actor-critic framework (Pinto et al., 2018; Andrychowicz et al., 2020; Baisero & Amato, 2021), which is also called asymmetric learning. This approach is particularly useful in multi-agent settings where global information is naturally accessible (Foerster et al., 2018; Lowe et al., 2017; Yu et al., 2022). Additional strategies include learning from noisy demonstrations (Tangkaratt et al., 2021), improving via self-correction from past trajectories (Song et al., 2018), and surpassing imperfect experts through regret-minimization frameworks (Chang et al., 2015). Recent work also explores leveraging LLMs as privileged experts, such as in embodied agents trained with reflective text-based guidance (Yang et al., 2024).

Besides, there are also representation learning techniques provided in order to reconstruct the privileged information (or its latent representation) via partial observation. For example Sermanet et al. (2018); Seo et al. (2023) use multi-view setups (e.g., image-based manipulation with additional camera views) to learn more informative embeddings. Others (Lee et al., 2020b; Salter et al., 2021; Kumar et al., 2021; Qi et al., 2023) leverage privileged simulator states during training and design policies that operate on both observed and inferred states.

In our experiments, we benchmark against several algorithms inspired by these lines of work, with detailed descriptions of the baselines provided in Appendix E.1.

B OMITTED PROOFS

Proposition 1. *If the guider’s policy is updated using policy mirror descent in each GPO iteration:*

$$\hat{\mu} = \arg \min \left\{ -\eta_k \langle \nabla V(\mu^{(k)}), \mu \rangle + \frac{1}{1-\gamma} D_{\mu^{(k)}}(\mu, \mu^{(k)}) \right\}, \quad (13)$$

then the learner’s policy update follows a constrained policy mirror descent:

$$\pi^{(k+1)} = \arg \min_{\pi \in \Pi} \left\{ -\eta_k \langle \nabla V(\pi^{(k)}), \pi \rangle + \frac{1}{1-\gamma} D_{\pi^{(k)}}(\pi, \pi^{(k)}) \right\} \quad (14)$$

Proof. First, since D is a weighted sum of KL divergence, it satisfies the definition of a Bregman divergence. Therefore, for any distributions $p, q \in \Delta(A)^{|S|}$, we have

$$D_q(p, q) = h_q(p) - h_q(q) - \langle \nabla h_q(q), p - q \rangle, \quad (15)$$

where $h_q(p) = \sum_{s \sim d_q} p_s \log p_s$ is the negative entropy weighted by the state distribution.

Next, by backtracking $\mu^{(k)}$ to $\pi^{(k)}$ from the last time step, we get:

$$\begin{aligned} \hat{\mu} &= \arg \min \left\{ -\eta_k \langle \nabla V(\mu^{(k)}), \mu \rangle + \frac{1}{1-\gamma} D_{\mu^{(k)}}(\mu, \mu^{(k)}) \right\} \\ &= \arg \min \left\{ -\eta_k \langle \nabla V(\pi^{(k)}), \mu \rangle + \frac{1}{1-\gamma} D_{\pi^{(k)}}(\mu, \pi^{(k)}) \right\} \\ &= \arg \min \left\{ -(1-\gamma)\eta_k \langle \nabla V(\pi^{(k)}), \pi \rangle + h_{\pi^{(k)}}(\pi) - \langle \nabla h_{\pi^{(k)}}(\pi^{(k)}), \pi \rangle \right\}, \end{aligned} \quad (16)$$

The optimality condition for $\hat{\mu}$ requires:

$$-(1-\gamma)\eta_k \nabla V(\mu^{(k)}) + \nabla h_{\mu^{(k)}}(\hat{\mu}) - \nabla h_{\mu^{(k)}}(\mu^{(k)}) = 0, \quad (17)$$

where we use the fact that:

$$\nabla_p D_q(p, q) = \nabla_p h_q(p) - \nabla_p h_q(q). \quad (18)$$

Now, consider the update of the learner’s policy, which involves a Bregman projection \mathcal{P}_Π :

$$\begin{aligned}
\pi^{(k+1)} &= \mathcal{P}_\Pi(\hat{\mu}) = \arg \min_{\pi \in \Pi} D_{\mu^{(k)}}(\pi, \hat{\mu}) \\
&= \arg \min_{\pi \in \Pi} \{h_{\mu^{(k)}}(\pi) - \langle \nabla h_{\mu^{(k)}}(\hat{\mu}), \pi \rangle\} \\
&= \arg \min_{\pi \in \Pi} \{h_{\pi^{(k)}}(\pi) - \langle \nabla h_{\pi^{(k)}}(\pi^{(k)}), \pi \rangle + (1 - \gamma)\eta_k \nabla V(\pi^{(k)}), \pi \rangle\} \\
&= \arg \min_{\pi \in \Pi} \{- (1 - \gamma)\eta_k \langle \nabla V(\pi^{(k)}), \pi \rangle + h_{\pi^{(k)}}(\pi) - \langle \nabla h_{\pi^{(k)}}(\pi^{(k)}), \pi \rangle\} \\
&= \arg \min_{\pi \in \Pi} \{-\eta_k \langle \nabla V(\pi^{(k)}), \pi \rangle + \frac{1}{1 - \gamma} D_{\pi^{(k)}}(\pi, \pi^{(k)})\}
\end{aligned} \tag{19}$$

This completes the proof. \square

Proposition 2. For policy π , μ , β and all state s , suppose $D_{TV}(\mu(\cdot|s), \beta(\cdot|s)) \lesssim \epsilon/2$, then we have

$$\mathbb{E}_{a \sim \beta} [1 - \rho^\pi(s, a)] \lesssim \epsilon + \sqrt{2d_{\text{targ}}}. \tag{20}$$

Proof. First, let’s examine the assumption $D_{TV}(\mu(\cdot|s), \beta(\cdot|s)) \lesssim \epsilon/2$ to check its validity.

Notice that at the start of each PPO policy update, the importance sampling ratio $\rho^\mu(s, a)$ equals 1 because the behavioral policy is equal to the policy being updated, i.e., $\beta(a|s) = \mu(a|s)$.

As PPO proceeds, $\rho^\mu(s, a)$ is updated multiple times using the same batch of samples. Due to the clipping function applied to $\rho^\mu(s, a)$, i.e., $\text{clip}(\rho^\mu(s, a), 1 - \epsilon, 1 + \epsilon)$, only state-action pairs for which $\rho^\mu(s, a) \in (1 - \epsilon, 1 + \epsilon)$ get updated. Hence, in the early epochs of PPO, with a properly tuned step size, we expect:

$$|1 - \rho^\mu(s, a)| \lesssim \epsilon. \tag{21}$$

Now, recalling the definition of total variation (TV) distance:

$$D_{TV}(\mu(\cdot|s), \beta(\cdot|s)) = \frac{1}{2} \sum_a |\mu(a|s) - \beta(a|s)| = \frac{1}{2} \sum_a \beta(a|s) |\rho^\mu(s, a) - 1| \lesssim \epsilon/2. \tag{22}$$

This confirms that the assumption $D_{TV}(\mu(\cdot|s), \beta(\cdot|s)) \lesssim \epsilon/2$ is reasonable, especially for the first few policy updates.

By the triangle inequality for total variation distance:

$$D_{TV}(\pi(\cdot|o), \beta(\cdot|s)) \leq D_{TV}(\pi(\cdot|o), \mu(\cdot|s)) + D_{TV}(\mu(\cdot|s), \beta(\cdot|s)), \tag{23}$$

we have

$$\begin{aligned}
D_{TV}(\pi(\cdot|o), \beta(\cdot|s)) &\leq \sqrt{\frac{1}{2} D_{\text{KL}}(\pi(\cdot|o), \mu(\cdot|s))} + D_{TV}(\mu(\cdot|s), \beta(\cdot|s)) \\
&\lesssim \sqrt{\frac{1}{2} d_{\text{targ}}} + \epsilon/2,
\end{aligned}$$

where we use Pinsker’s inequality to bound the total variation distance between π and μ in terms of their KL divergence.

Finally, since total variation is linked to the expected difference between probabilities under different policies, we have:

$$\mathbb{E}_{a \sim \beta} [1 - \rho^\pi(s, a)] = 2D_{TV}(\pi(\cdot|o), \beta(\cdot|s)) \lesssim \epsilon + \sqrt{2d_{\text{targ}}}. \tag{24}$$

This result implies that, under the assumption, the majority of samples are valid for updating the learner’s policy during the early PPO epochs. \square

C PSEUDO CODE

In this section, we present the pseudo code of our algorithm (see Algorithm 1). The algorithm is based on PPO, with an additional objective to leverage the extra information available during training.

Algorithm 1: Guided Policy Optimization

Input: Initial policy parameters θ_0 , value function parameters ϕ_0
for $k = 0, 1, 2, \dots$ **do**
 Collect trajectory set $\mathcal{D}_K = \{\tau_i\}$ by running guider policy $\mu_k = \mu(\cdot|o_g; \theta_k)$;
 Compute rewards-to-go \hat{R}_t ;
 Compute advantage estimates \hat{A}_t using GAE w.r.t. current value function V_{ϕ_k} ;
 Update policy parameters θ_k to θ_{k+1} by maximizing GPO objective (11) or (12);
 Fit value function parameters ϕ_{k+1} by minimizing mean squared error;

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_K|T} \sum_{\tau \in \mathcal{D}_K} \sum_{t=0}^T \left(V_{\phi}((o_g)_t) - \hat{R}_t \right)^2$$

Table 4: TigerDoor-alt problem

state \ action	action	
	a_L	a_R
s_L	2	0
s_R	0	1

D GPO ON TIGERDOOR-ALT PROBLEM

We introduce an alternative version of the TigerDoor, called TigerDoor-alt, which also highlights an imitation gap, even without additional exploratory information. In this scenario, the listening action a_l is removed, and the reward for correctly selecting the left door is increased to 2 as shown in Table 3. Similarly, the teacher continues to select the correct door, while the student learns to randomly choose between the two doors, yielding an expected reward of 0.75. However, the optimal policy for the student is to always choose the left door, which provides an expected reward of 1. This discrepancy arises from the loss of information when converting the reward-based objective into a policy-supervised objective.

Now we provide an intuitive example to illustrate how GPO can achieve the optimal policy in the TigerDoor-alt problem. At time step t , suppose the guider’s and learner’s policies are:

$$\mu_t(\cdot|s_L) = \mu_t(\cdot|s_R) = \pi_t = (x_t, y_t),$$

where the two policies are equal due to the backtracking from time step $t - 1$. After one update step, the guider’s policy becomes:

$$\hat{\mu}_t(\cdot|s_L) = (x_t + p_t, y_t - p_t), \quad \hat{\mu}_t(\cdot|s_R) = (x_t + q_t, y_t - q_t)$$

The key insight is that the higher reward for (s_L, a_L) compared to (s_R, a_R) leads to a larger gradient step, implying $p_t > q_t$. The learner then imitates the guider, resulting in the updated policy:

$$\pi_{t+1} = \left(x_t + \frac{p_t - q_t}{2}, y_t - \frac{p_t - q_t}{2} \right).$$

Hence, $\pi_{t+1}(a_L) > \pi_t(a_L)$, meaning the learner’s policy moves closer to the optimal policy $(1, 0)$ in a monotonic fashion.

The critical mechanism is that actions yielding higher rewards induce larger updates in the guider’s policy, which the learner then captures via imitation. Meanwhile, the backtracking step keeps the guider aligned with the learner, enabling steady and consistent policy improvement.

E EXPERIMENTAL SETTINGS**E.1 BASELINES**

In this section, we briefly introduce the baselines used in our experiments.

PPO. This is the standard algorithm used to train the agent without any additional information. The objective function is given by:

$$\mathcal{L}(\pi) = -\mathbb{E} \left[\min \left(\rho^\pi(o_l, a) A^\beta(o_l, a), \rho_{clip}^\pi(o_l, a, \epsilon) A^\beta(o_l, a) \right) \right], \quad (25)$$

where the behavioral policy is $\beta = \pi_{old}$.

GPO-naive. This variant of GPO uses the GPO-penalty without the auxiliary RL loss term. The objective function is:

$$\begin{aligned} \mathcal{L}_{\text{GPO-naive}}(\theta) = -\mathbb{E} \left[\min \left(\rho^{\mu_\theta} A^\beta(o_g, a), \rho_{clip}^{\mu_\theta} A^\beta(o_g, a) \right) - \alpha \text{D}_{\text{KL}}(\mu_\theta(\cdot|o_l) || \pi_{\hat{\theta}}(\cdot|o_g)) \right. \\ \left. - \text{D}_{\text{KL}}(\mu_{\hat{\theta}}(\cdot|o_l) || \pi_\theta(\cdot|o_g)) \right]. \end{aligned} \quad (26)$$

GPO-ablation. This is another variant of GPO-penalty, but without the BC loss term. The objective is:

$$\begin{aligned} \mathcal{L}_{\text{GPO-ablation}}(\theta) = -\mathbb{E} \left[\min \left(\rho^{\mu_\theta} A^\beta(o_g, a), \rho_{clip}^{\mu_\theta} A^\beta(o_g, a) \right) - \alpha \text{D}_{\text{KL}}(\mu_\theta(\cdot|o_l) || \pi_{\hat{\theta}}(\cdot|o_g)) \right. \\ \left. + \min \left(\rho^{\pi_\theta} A^\beta(o_g, a), \rho_{clip}^{\pi_\theta} A^\beta(o_g, a) \right) \right]. \end{aligned} \quad (27)$$

PPO-asym. This method trains the student using PPO, but with asymmetric value function taking o_g as input. The objective is:

$$\mathcal{L}(\pi) = -\mathbb{E} \left[\min \left(\rho^\pi(o_l, a) A^\beta(o_g, a), \rho_{clip}^\pi(o_l, a, \epsilon) A^\beta(o_g, a) \right) \right]. \quad (28)$$

ADVISOR. Given teacher’s policy μ , ADVISOR (Weihs et al., 2024) uses a balancing coefficient w between BC and RL training, based on the distance between the teacher’s policy μ and an auxiliary imitation policy $\hat{\pi}$:

$$\mathcal{L}(\pi) = -\mathbb{E} \left[w \text{CE}(\mu(\cdot|o_g), \pi(\cdot|o_l)) + (1 - w) \min \left(\rho^\pi(o_l, a) A^\beta(o_l, a), \rho_{clip}^\pi(o_l, a, \epsilon) A^\beta(o_l, a) \right) \right],$$

where $w = \exp(-\alpha \text{D}_{\text{KL}}(\mu(\cdot|o_g), \hat{\pi}(\cdot|o_l)))$ and CE means cross-entropy.

ADVISOR-co. This is a modified version of the ADVISOR algorithm for co-training setting, as the original does not involve teacher training. The teacher’s objective is:

$$\mathcal{L}(\mu) = -\mathbb{E} \left[\min \left(\rho^\mu(o_g, a) A^\beta(o_g, a), \rho_{clip}^\mu(o_g, a, \epsilon) A^\beta(o_g, a) \right) \right]. \quad (29)$$

ADVISOR-co can be viewed as GPO-penalty without the backtrack term and with a different α -update schedule. However, in the absence of backtracking, the coefficient w quickly diminishes, as the auxiliary policy cannot follow the teacher effectively, reducing this approach to pure PPO training for the student.

PPO+BC. In this approach, the teacher is trained using PPO:

$$\mathcal{L}(\mu) = -\mathbb{E} \left[\min \left(\rho^\mu(o_g, a) A^\beta(o_g, a), \rho_{clip}^\mu(o_g, a, \epsilon) A^\beta(o_g, a) \right) \right], \quad (30)$$

while the student is trained using BC with the teacher:

$$L(\pi) = \mathbb{E} [\text{D}_{\text{KL}}(\mu(\cdot|o_g), \pi(\cdot|o_l))]. \quad (31)$$

PPO+BC-t. Given teacher’s policy μ , the student is trained using a combined loss of PPO and BC:

$$\mathcal{L}(\pi) = -\mathbb{E} \left[\min \left(\rho^\mu(o_g, a) A^\beta(o_g, a), \rho_{clip}^\mu(o_g, a, \epsilon) A^\beta(o_g, a) \right) - \text{D}_{\text{KL}}(\mu(\cdot|o_g), \pi(\cdot|o_l)) \right]. \quad (32)$$

A2D. Adaptive Asymmetric DAgger (A2D) (Warrington et al., 2020) is closely related to GPO, as it also involves co-training both the teacher and the student. A2D uses a mixture policy $\beta(a|o_g, o_l) =$

$\lambda\mu(a|o_g) + (1 - \lambda)\pi(a|o_l)$ to collect trajectories and train the expert μ with a mixed value function $V(o_g, o_l) = \lambda V^\mu(o_g) + (1 - \lambda)V^\pi(o_l)$. The objective is:

$$\mathcal{L}(\mu) = -\mathbb{E} \left[\min \left(\rho^\mu(o_g, o_l, a) A^\beta(o_g, o_l, a), \rho_{clip}^\mu(o_g, o_l, a, \epsilon) A^\beta(o_g, o_l, a) \right) \right], \quad (33)$$

while the student is updated through BC:

$$\mathcal{L}(\pi) = \mathbb{E} [\text{D}_{\text{KL}}(\mu(\cdot|o_g), \pi(\cdot|o_l))] \quad (34)$$

In practice, A2D sets $\lambda = 0$ or anneals it quickly for better performance. When $\lambda = 0$, A2D is equivalent to GPO-naive without the backtrack step, and it uses the student’s behavioral policy π instead of the teacher’s policy μ . While A2D implicitly constrains the teacher’s policy through the PPO clipping mechanism (which prevents the teacher from deviating too far from the student’s behavioral policy), this is insufficient to replace the explicit backtrack step. The gap between μ and π can accumulate if the student fails to follow the teacher. Consequently, most samples will be clipped as training progresses, leading A2D to struggle in training a strong teacher.

ELF. Given teacher’s policy μ , ELF Distillation (Walsman et al., 2023) trains two policies jointly: a follower π_f to mimic the teacher through BC:

$$\mathcal{L}(\pi_f) = \mathbb{E} [\text{D}_{\text{KL}}(\mu(\cdot|o_g), \pi_f(\cdot|o_l))], \quad (35)$$

and an explorer π_e trained through PPO:

$$\mathcal{L}(\pi_e) = -\mathbb{E} \left[\min \left(\rho^{\pi_e}(o_l, a) A(o_l, a), \rho_{clip}^{\pi_e}(o_l, a, \epsilon) A(o_l, a) \right) \right], \quad (36)$$

To utilize teacher supervision, ELF applies a potential-based reward shaping (Ng et al., 1999) $r + \gamma V^{\pi_f}(o'_l) - V^{\pi_f}(o_l)$ to the explorer, where V^{π_f} is the value function of follower. However, ELF needs to divide the interaction equally to train the follower and the explorer, which leads to inefficiencies.

ELF-asym. Since the follower is not required during execution, an asymmetric value function $V^{\pi_f}(o_g)$ is used instead of the original one. Although there are some performance improvement, ELF-asym still performs worse than PPO-asym due to inefficient experience usage.

L2T-PPO. Similar to PPO+BC, the teacher is trained using PPO:

$$\mathcal{L}(\mu) = -\mathbb{E} \left[\min \left(\rho^\mu(o_g, a) A^\beta(o_g, a), \rho_{clip}^\mu(o_g, a, \epsilon) A^\beta(o_g, a) \right) \right], \quad (37)$$

while the student is trained using a combined loss of PPO and BC with the teacher:

$$\mathcal{L}(\pi) = -\mathbb{E} \left[\min \left(\rho^\mu(o_g, a) A^\beta(o_g, a), \rho_{clip}^\mu(o_g, a, \epsilon) A^\beta(o_g, a) \right) - \text{D}_{\text{KL}}(\mu(\cdot|o_g), \pi(\cdot|o_l)) \right], \quad (38)$$

where the behavioral policy $\beta = \mu$.

E.2 HYPERPARAMETERS

The experiments in Sections 4.1 and 4.3 use the same codebase from (Lu et al., 2023). The hyperparameters for these experiments are listed in Table 5. For GPO-clip, due to the asymmetry with large δ , we replace the $\text{clip}(\frac{\mu}{\pi}, 1 - \delta, 1 + \delta)$ with $\text{clip}(\frac{\mu}{\pi}, \frac{1}{\rho}, \rho)$ in the POPGym tasks.

For the experiments in Section 4.2, we use the codebase from (Freeman et al., 2021). We perform a hyperparameter search for the original versions of the tasks and then fix the same hyperparameters for the partially observable and noisy variants. The hyperparameter search is detailed in Table 6, and the selected hyperparameters for the experiments are provided in Table 7. Other fixed hyperparameters are listed in Table 8.

All algorithms in Brax are run with 10 random seeds, whereas those in POPGym use 3 seeds, as the latter exhibits lower variance. Reward curves in this paper report the mean and standard deviation across these runs.

Table 5: Hyperparameters used in TigerDoor and POPGym.

Parameter	Value (TigerDoor)	Value (POPGym)
Adam Learning Rate	5e-5	5e-5
Number of Environments	64	64
Unroll Length	1024	1024
Number of Timesteps	2e6	15e6
Number of Epochs	30	30
Number of Minibatches	8	8
Discount γ	0.99	0.99
GAE λ	1.0	1.0
Clipping Coefficient ϵ	0.2	0.2
Entropy Coefficient	0.0	0.0
Value Function Weight	1.0	1.0
Maximum Gradient Norm	0.5	0.5
Activation Function	LeakyReLU	LeakyReLU
Encoder Layer Sizes	128	[128,256]
Recurrent Layer Hidden Size	-	256
Action Decoder Layer Sizes	128	[128,128]
Value Decoder Layer Sizes	128	[128,128]
KL Threshold d	0.001	0.1 (0.001 for CartPole)
Clip ρ	1.1	10 (1.2 for CartPole)
RL Coefficient α	1	0 (1 for CartPole)

Table 6: Sweeping procedure in the Brax domain.

Parameter	Value
Reward Scaling r_s	[0.1, 1]
Discount γ	[0.97, 0.99, 0.997]
Unroll Length l	[5, 10, 20]
Batchsize b	[256, 512, 1024]
Number of Minibatches n	[4, 8, 16, 32]
Number of Epochs e	[2, 4, 8]
Entropy Coefficient c	[0.01, 0.001]
KL Threshold d	[0.01, 0.001]
Clip δ	[0.1, 0.3]
RL Coefficient α	[0, 2, 3]

E.3 ENVIRONMENT DESCRIPTIONS

We provide a brief overview of the environments used and the guider’s observation settings.

Brax tasks and CartPole in POPGym: For these tasks, velocities and angular velocities are removed from the learner’s observation. Gaussian noise with standard deviations of 0.1, 0.2, and 0.3 is added to the observations, corresponding to the difficulty levels *Easy*, *Medium*, and *Hard*, respectively. The guider, however, has access to the noiseless observations and the removed velocities.

Autoencode: During the WATCH phase, a deck of cards is shuffled and played in sequence to the agent with the watch indicator set. The watch indicator is unset at the last card in the sequence, where the agent must then output the sequence of cards in order. The guider directly observes the correct card to be output at each timestep.

Battleship: A partially observable version of Battleship game, where the agent has no access to the board and must derive its own internal representation. Observations contain either HIT or MISS and the position of the last salvo fired. The player receives a positive reward for striking a ship, zero reward for hitting water, and negative reward for firing on a specific tile more than once. The guider has access to a recorder that tracks all previous actions taken by the agent.

Count Recall: Each turn, the agent receives a next value and query value. The agent must answer the query with the number of occurrences of a specific value. In other words, the agent must store

Table 7: Adopted hyperparameters in the Brax domain. Notations correspond to Table 6.

Task	r_s	γ	l	b	n	e	c	d	δ	α
Ant	0.1	0.97	5	1024	32	4	0.01	0.001	0.3	2
Halfcheetah	1	0.99	5	512	4	4	0.001	0.001	0.1	2
Humanoid	0.1	0.99	5	512	32	4	0.01	0.001	0.1	2
HumanoidStandup	0.1	0.99	5	256	32	8	0.01	0.001	0.3	3
InvertedDoublePendulum	1	0.997	20	256	8	4	0.01	0.001	0.1	0
Swimmer	1	0.997	5	256	32	4	0.01	0.001	0.3	3
Walker2d	1	0.99	5	512	32	4	0.01	0.001	0.1	2

Table 8: Common hyperparameters used in Brax domain.

Parameter	Value
Adam Learning Rate	3e-4
Number of Environments	2048
Episode Length	1024
Number of Timesteps	3e7
GAE λ	0.95
Clipping Coefficient ϵ	0.3
Activation Function	SiLU
Value Layer Sizes	[128, 128]
Policy Layer Sizes	[128, 128]

running counts of each unique observed value, and report a specific count back, based on the query value. The guider directly observes the running counts at each timestep.

Repeat Previous: At the first timestep, the agent receives one of four values and a remember indicator. Then it randomly receives one of the four values at each successive timestep without the remember indicator. The agent is rewarded for outputting the observation from some constant k timesteps ago, i.e. observation o_{t-k} at time t . The guider has direct access to the value o_{t-k} at time t .

E.4 ADDITIONAL COMPARATIVE EXPERIMENTS

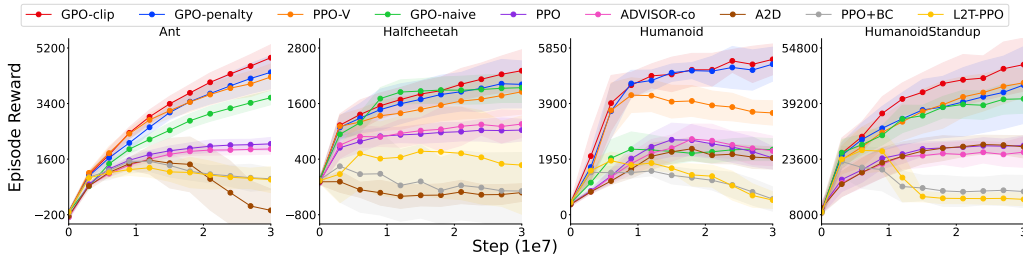


Figure 8: Performance comparison on selected Brax tasks, including L2T-PPO.

Here, we present additional baselines that were not included in the main paper. First, L2T-RL tackles a problem similar to ours by employing a fully observable teacher to supervise a partially observable student. However, L2T-RL resembles PPO+BC, as it uses the teacher purely as a behavioral policy without aligning it with the student’s policy, which results in ineffective supervision. Moreover, L2T-RL relies on the teacher’s experience to train the student through RL without any offline adaptation, further limiting the effectiveness of the RL training. Fig. 8 illustrates the performance of the PPO-based L2T-RL (details in Appendix E.1), where L2T-PPO performs similarly to PPO+BC and falls short of the other methods proposed in our work.

Second, we provide a more detailed comparison with methods that train a teacher first and then apply TSL techniques to address the challenge of an inimitable teacher. Two state-of-the-art approaches in this category are ADVISOR and TGRL, which are based on PPO and SAC, respectively. We evaluate these methods using a PPO-trained teacher with full observability on the Ant task. The results, shown

in Fig. 9, indicate that while both methods perform well and demonstrate improved efficiency over their respective base algorithms under full observability, their performance degrades in the partially observable Ant task. In this case—where the teacher’s policy is effectively inimitable—the TSL methods perform comparably to their base RL algorithms.

Fig. 10 shows the KL divergence between the agent policies of these TSL methods and the teacher. Under full observability, where the teacher was trained, the agents can successfully mimic the teacher’s policy. However, under partial observability, the agents struggle to imitate the teacher’s behavior, leading to a substantial KL divergence. Since both ADVISOR and TGRL revert to standard RL when the teacher becomes inimitable, this explains their performance similarity to the base algorithms in such scenarios.

Additionally, we report TGRL’s performance across the 28 Brax tasks used in our experiments (see Fig. 11). Note that TGRL follows an off-policy training paradigm, in contrast to all other methods presented in the main paper, which makes it significantly slower (refer to Table 9). Therefore, we run TGRL for only 1M steps, which is sufficient for convergence as shown by the learning curves.

Finally, we include the representation learning method RMA (Kumar et al., 2021), which aims to reconstruct the latent privileged information used by the teacher during the student training phase. Such representation learning approaches are promising when privileged information can be reliably approximated from partial observations. However, their effectiveness is limited in Brax, where observations are noisy (Figure 11). Since the noise cannot be removed by a deterministic mapping without additional structure or assumptions, regression-based reconstruction tends to collapse to identity mappings and fails to recover meaningful latent representations. Moreover, current encoder-based pipelines often lack theoretical guarantees for convergence or generalization across diverse tasks, particularly under partial observability.

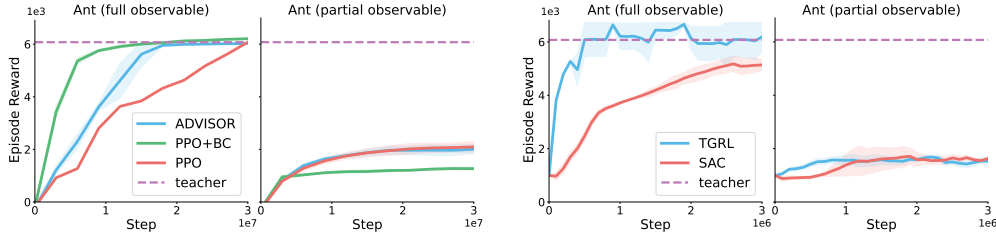


Figure 9: ADVISOR, PPO+BC and TGRL with a pre-trained teacher.

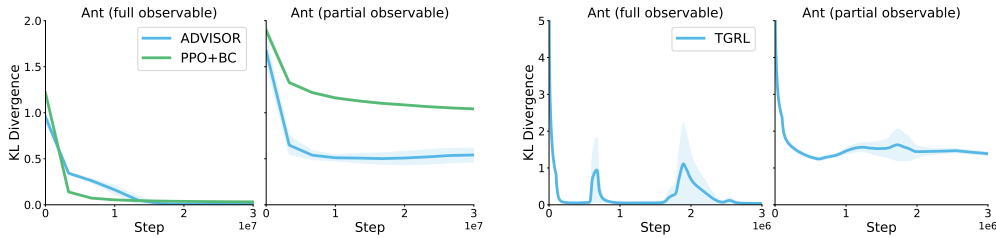


Figure 10: The KL divergence of ADVISOR, PPO+BC and TGRL with pre-trained teacher.

E.5 ADDITIONAL FIGURES

Fig. 12 shows the reward curves of the experiments presented in Section 4.2. Fig. 13 illustrates the performance influenced by the parameter sharing. We can observe that parameter sharing can sometimes impair performance, particularly when the observation dimension is large. For instance, in the *HumanoidStandup* task, the observation dimension is 400, which challenges the expressive capacity of the network. Thus, the decision to share the policy network represents a trade-off between memory efficiency and performance.

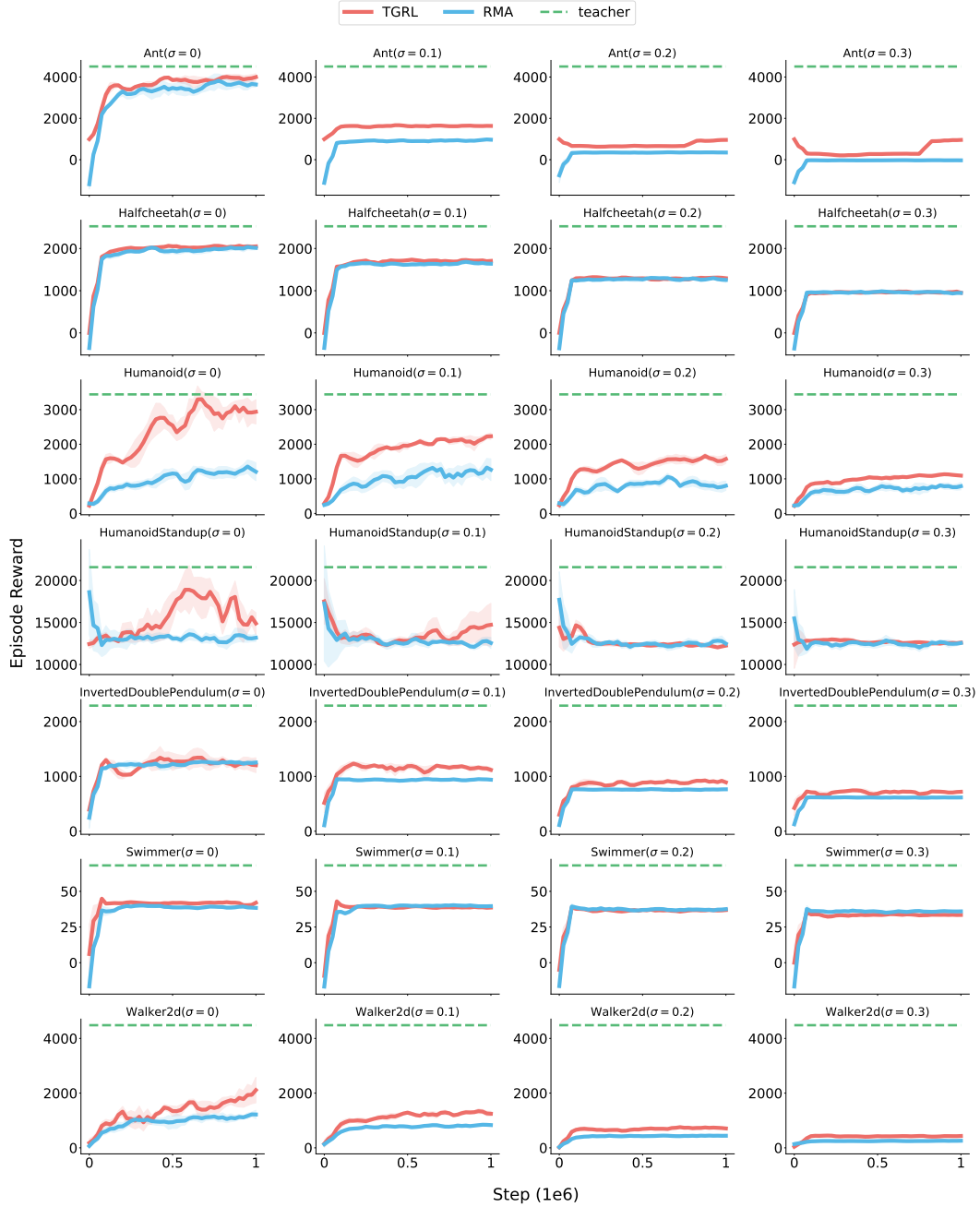


Figure 11: TGRL and RMA with pre-trained teacher on 28 Brax tasks.

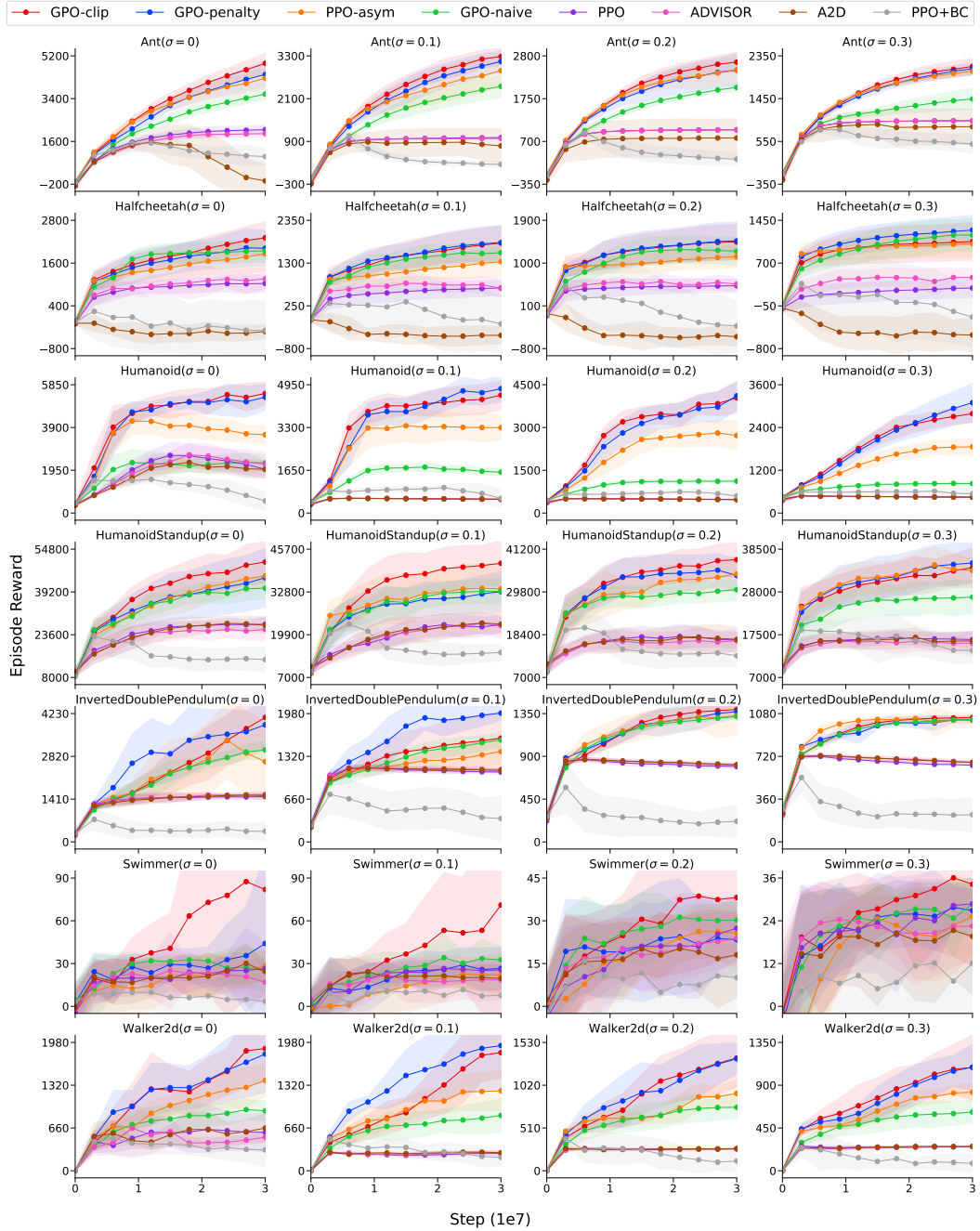


Figure 12: Comparing between GPO and other baselines on 28 Brax tasks.

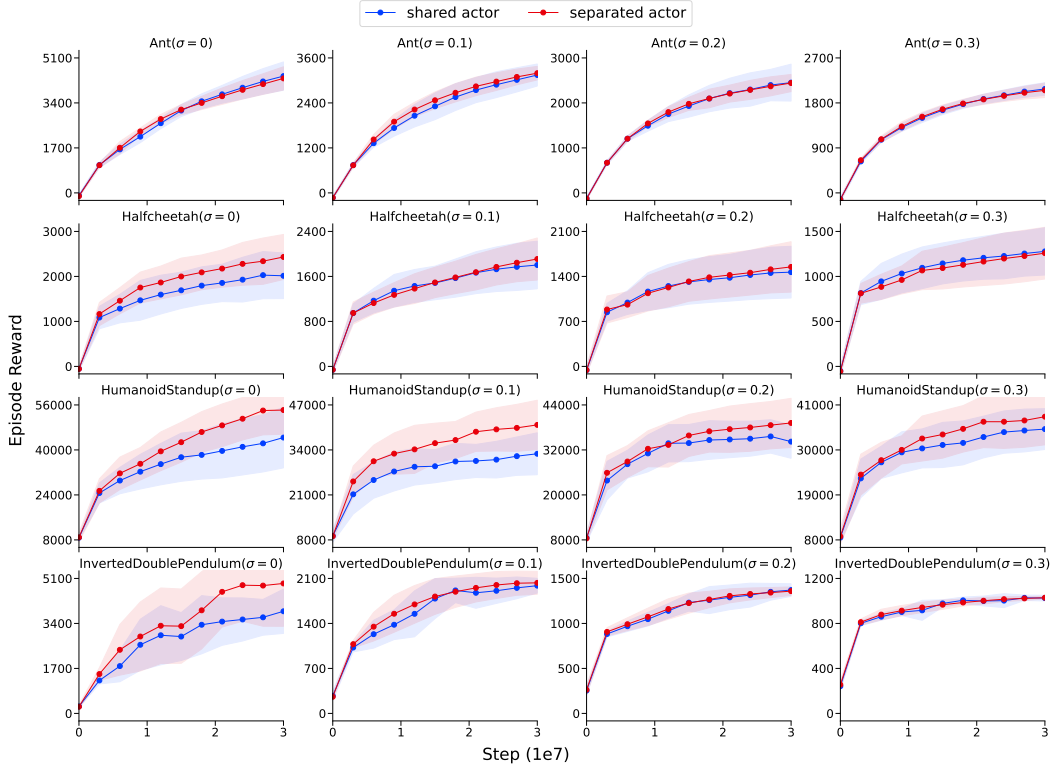


Figure 13: Comparing shared and separated policy networks of GPO-penalty.

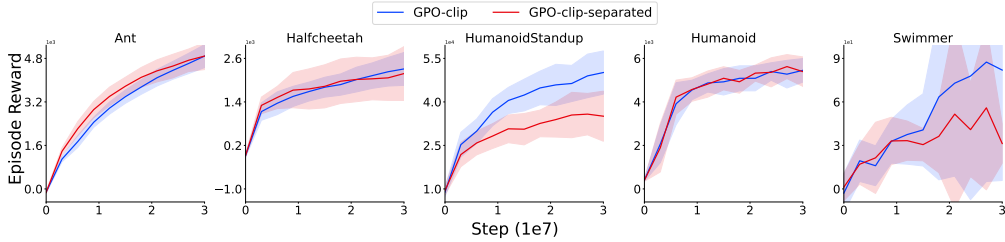


Figure 14: Comparing joint and separated update of GPO-clip.

E.6 COMPUTATIONAL COST

In this section, we compare the computational cost of GPO (both GPO-penalty and GPO-clip share the same cost), PPO-asym, TGRL and the environmental step time across several environments. The results, presented in Table 9, show that GPO is approximately 10% to 20% slower than PPO-asym. Importantly, GPO achieves this with no additional networks, underscoring its efficiency despite the modest increase in computational overhead.

Table 9: Frames Per Second (FPS) of GPO, PPO-asym and TGRL across several environments, computed on the NVIDIA GeForce RTX 4090.

Environment	GPO	PPO-asym	TGRL	Environmental Step
Ant	1.19×10^5	1.36×10^5	1.13×10^2	4.23×10^5
Halfcheetah	6.27×10^4	7.21×10^4	9.58×10^1	2.55×10^5
Humanoid	6.29×10^4	7.18×10^4	9.92×10^1	2.50×10^5
Swimmer	3.33×10^4	3.83×10^4	1.04×10^2	1.50×10^5

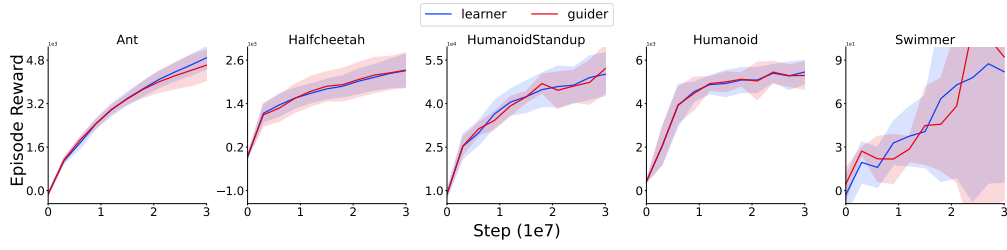


Figure 15: Comparing the performance of guider and learner of GPO-clip.