

# SHIRO: Soft Hierarchical Reinforcement Learning

**Kandai Watanabe**

Department of Computer Sciences  
University of Colorado Boulder  
United States  
kandai.watanabe@colorado.edu

**Mathew Strong**

Department of Computer Sciences  
University of Colorado Boulder  
United States  
mathew.strong@colorado.edu

**Omer Eldar**

Department of Computer Sciences  
University of Colorado Boulder  
United States  
omer.eldar@colorado.edu

## Abstract:

Hierarchical Reinforcement Learning (HRL) algorithms have been demonstrated to perform well on high-dimensional decision making and robotic control tasks. However, because they solely optimize for rewards, the agent tends to search the same space redundantly. This problem reduces the speed of learning and achieved reward. In this work, we present an Off-Policy HRL algorithm that maximizes entropy for efficient exploration. The algorithm learns a temporally abstracted low-level policy and is able to explore broadly through the addition of entropy to the high-level. The novelty of this work is the theoretical motivation of adding entropy to the RL objective in the HRL setting. We empirically show that the entropy can be added to both levels if the Kullback-Leibler (KL) divergence between consecutive updates of the low-level policy is sufficiently small. We performed an ablation study to analyze the effects of entropy on hierarchy, in which adding entropy to high-level emerged as the most desirable configuration. Furthermore, a higher temperature in the low-level leads to Q-value overestimation and increases the stochasticity of the environment that the high-level operates on, making learning more challenging. Our method, SHIRO, surpasses state-of-the-art performance on a range of simulated robotic control benchmark tasks and requires minimal tuning.

**Keywords:** Off-Policy, Hierarchical, Reinforcement Learning, Maximum Entropy

## 1 Introduction

Deep Reinforcement Learning (RL) is a promising approach for robots to learn complex policies; e.g., manipulators [1], quadrupedal robots [2, 3] and humanoids [4, 5, 6]. Learning complex policies is made possible with the help of strong function approximators such as neural networks. This enables an agent to learn a policy in both a high dimensional state and action space [7]. Deep RL methods are well-suited for learning one atomic skill; e.g., a manipulator reaching an object or a quadruped walking upright. While very powerful, such methods struggle to learn compositional, long-horizon skills, such as utilizing locomotion for solving a maze. To solve such tasks, a better structure and stronger exploration strategies are required.

To this end, Hierarchical Reinforcement Learning (HRL) has been empirically shown to work well on these tasks by making use of the temporal abstraction generated from its hierarchical structure [8, 9, 10, 11]. In HRL, policies are stacked into multiple levels to construct a hierarchy, where the low-level agent receives information from the one(s) above it. This allows for the decomposition of a long-horizon task into smaller, more easily solvable problems. The theory of HRL has been well-studied over the years, but has lacked a generalized method that can train an agent efficiently in a complex, continuous space [12, 13, 14, 15, 16].

One potential solution is the state-of-the-art algorithm named Hierarchical Reinforcement Learning With Off-Policy Correction (HIRO) [8]. It proposes to stack two off-policy Deep RL agents and to relabel past experiences in the high-level policy’s replay buffer to correct the changes in the low-level policy. This allows training both levels simultaneously using an off-policy training method, which speeds up the training. Accordingly, this method has been extended to new algorithms [11, 10, 3, 17, 18]. However, because the agent solely optimizes for reward, it over-commits to actions it believes are optimal, leading to poor exploration. This leads to inefficient training and slower convergence.

We present an off-policy hierarchical reinforcement learning algorithm (SHIRO: Soft Hierarchical Reinforcement Learning with Off-Policy Correction). We propose adding entropy maximization to the objective [19] of both levels of the agent to improve the training efficiency of HIRO. The addition of entropy in the RL objective has demonstrated an improvement of efficiency in exploration [20, 21, 22, 23, 24, 19, 25]. Additionally, we show that a relatively small Kullback-Leibler (KL) divergence in the low-level agent over time guarantees that the low-level agent can be viewed as part of the environment. This allows us to train the high-level agent as if it was a single agent and applies the soft policy iteration theorem in SAC to our work. We empirically demonstrate more efficient learning with the addition of entropy to the RL objective. Our method produces superior results in data efficiency and training time to the base implementation of HIRO on HRL benchmark environments for robotic locomotion (AntMaze, AntFall, and AntPush). Finally, we conduct an ablative study, where we apply different combinations of entropy, and experiment with making the entropy temperature a learned parameter. The results of the ablative study are used to produce the best possible method.

## 2 Related Work

The problem of finding a suitable hierarchical structure to solve a long-horizon task has been widely explored in the RL community for many years [12, 13, 14, 15]. General structures are a tree-like structure where multiple policies are stacked like a tree and a tower-like structure where single policies are stacked on top of each other. Classic tree-like works train a high-level policy, given a set of hand-engineered low-level policies [26, 27, 28]. One popular framework in a tree-like structure is the Options framework [14, 29]. This algorithm learns to pick an option and follows the selected option policy until it meets its termination condition. We can view the high-level policy as a switch over options. However, it requires prior knowledge for designing options. To this end, the option-critic [30] proposes a method to train the high-level policy jointly. Recent work on the Options framework utilizes regularizers to learn multiple sub-policies [30, 31, 32].

Recent works focus on designing rewards to learn sub-policies more effectively. One choice is to learn both high-level and low-level policies from final task rewards end-to-end [33, 34], including the aforementioned option-critic frameworks [30, 31]. Another major stream is to provide auxiliary rewards to low-level policies to foster the learning [28, 35, 36, 37, 38], such as hand-engineered low-level reward based on domain knowledge [39, 35, 36, 37], mutual information for more efficient exploration [40, 38] and goal-conditioned rewards [12, 41, 42, 9, 43, 8, 11].

The tower-like structure, such as the FeUdal framework [12, 43] and HIRO algorithms [8], generally takes a form of goal-conditioned rewards, so that the high-level policy acts as a Manager/Planner to lead the low-level policy to work for/reach the sub-goal that the high-level policy provided. The benefit of these methods is that the representation is generic and does not require hand-engineering specific to each domain. Especially, HIRO [8] uses the state in its raw form to construct a parameterized reward and [11, 10] extend the sub-goal to latent spaces. This goal-conditioned framework is also extended to the real robotic domain [3, 44] and extensively explored outside of HRL [45, 46, 41, 42, 47]. Problems of these algorithms are also targeted as in [48, 49, 50, 51], indicating that there are numerous avenues for improvement.

Furthermore, while very powerful, HIRO searches the same states redundantly, leading to inefficiency. Our algorithm extends the work to maximize entropy in addition to the original RL objective, that is shown to work well in other works [20, 21, 22, 23, 24, 19, 25], including finding diverse goals to speed up learning [52]. It is also applied in HRL [53] and also in the goal-conditioned setting [54, 55]. While [54, 55] are similar to our work, the intrinsic reward for the low-level controller must be hand-engineered and the controller is pre-trained to obtain enough abstraction before the joint

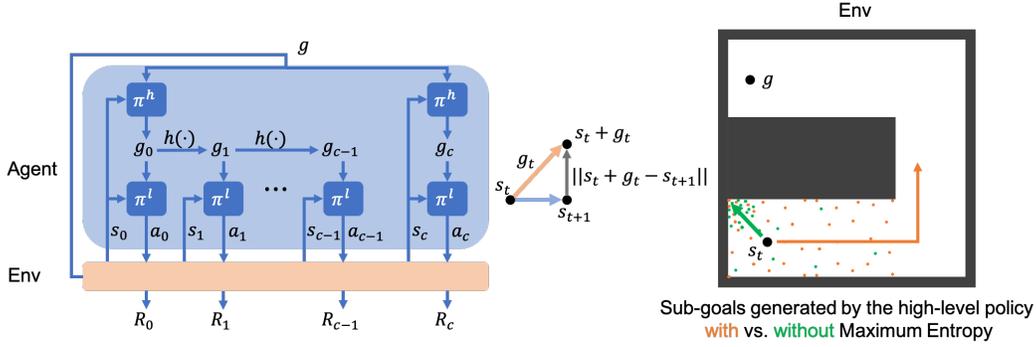


Figure 1: The left figure is a schematic of how the agent interacts with the environment. The high-level policy generates a sub-goal to guide the low-level policy in its direction as shown in the center figure. The right figure is the depiction of the high-level policy generating sub-goals at time  $t$  in AntMaze (from Figure 2). The entropy maximization helps the policies to diversify their sub-goals/actions for better exploration, allowing to get out of the local optima, faster convergence, and improved training of the low-level policy.

training. [18] is the closest to our work where they use SAC on both levels, but they don't explain on what condition entropy works well. While maximizing entropy is effective, too high entropy leads to poor performance. Regarding the problems of sub-goals, [17] added a constraint to avoid from generating unreachable sub-goals, and [52] minimizes changes in the high-level policy which facilitates learning. Unlike other methods, our work adds entropy to either or both levels and explains on what condition adding entropy works and how to best utilize entropy.

### 3 Preliminaries

#### 3.1 Goal Conditioned Reinforcement Learning

Goal-Conditioned Reinforcement Learning (RL) extends the definition of the RL problem by conditioning its policy and reward function on goals. To define formally, a Markov Decision Process (MDP) is augmented with a set of goals  $\mathcal{G}$ , and is called a *Universal MDP* (UMDP). A UMDP is a tuple containing  $(\mathcal{S}, \mathcal{G}, \mathcal{A}, P_T, R)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $P_T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability and  $R : \mathcal{S} \times \mathcal{G} \times \mathcal{A} \rightarrow [r_{min}, r_{max}]$  is the reward function mapping state, goal and action pair to a scalar value reward between  $r_{min}$  and  $r_{max}$ . The reward function is conditioned on a goal to represent how well the agent is approaching the goal. To approach each goal in a different manner, the agent's policy  $a_t = \pi(a_t | s_t, g)$  is augmented to condition on the goal. At the start of every episode, a goal  $g$  is sampled from a set of goals  $\mathcal{G}$  by the distribution  $P_g(g)$ . At every time step  $t$ , the environment generates state  $s_t \in \mathcal{S}$ . Then, the agent takes action  $a_t \in \mathcal{A}$  sampled from its policy  $\pi(a_t | s_t, g)$ , transitions to the next state  $s_{t+1}$  with probability  $P_T(s_{t+1} | s_t, a_t)$ , and receives a reward  $R_t = R(s_t, a_t, g)$ . The probability distribution over the samples generated by the policy  $\pi$  is denoted as  $\rho^\pi$ . The objective is to find the optimal policy  $\pi_\phi$  with parameters  $\phi$  that maximizes the expected reward  $\mathbb{E}_{(s_t, a_t) \sim \rho^\pi, g \sim P_g} [\sum_t \gamma^t R(s_t, a_t, g)]$ .

#### 3.2 Goal Conditioned Hierarchical Reinforcement Learning

HRL is an RL approach that leverages the hierarchical structure in an agent's model to decompose the original problem into sub-problems. This allows the algorithm to learn within the space of each sub-problem, and to learn a policy that can be shared across the sub-problems. Hierarchical Reinforcement learning with Off-policy correction (HIRO) [8] is a potential solution to this problem that makes use of goal-conditioned policies. This algorithm successfully accomplishes continuous complex tasks and is widely used [10, 56, 11]. Our work is based on their work, and is depicted in Figure 1. In the HIRO algorithm, the agent consists of a high-level policy  $\pi^h(g_t | s_t, g)$  and a low-level policy  $\pi^l(a_t | s_t, g_t)$ , where  $g_t$  is a sub-goal sampled from the high-level policy. At every  $c$  steps, the high-level policy computes a sub-goal for the low-level policy given the current state  $s_t$  and final goal  $g$ . The low-level policy is conditioned on the current state  $s_t$  and the sub-goal  $g_t$  to

generate an action  $a_t$ . For the rest of the steps, the sub-goal is advanced by the goal transition model  $h(s_t, g_t, s_{t+1}) = s_t + g_t - s_{t+1}$  based on the current and next state. The intrinsic reward is defined as  $r(s_t, g_t, a_t, s_{t+1}) = -\|s_t + g_t - s_{t+1}\|_2$  to reward the low-level policy, while the environment reward is for the high-level policy. To train the policies, the experiences are stored into replay buffers and apply off-policy training individually. HIRO addresses the non-stationary problem by relabeling past replay experiences with new goals that would prompt the current low-level policy to behave in the same way as the low-level policy in the past.

### 3.3 Soft Actor Critic

Recent advancements in off-policy RL have enabled policies to be trained on a range of challenging domains, from games to robotic control. Off-policy algorithms such as Q-learning, DDPG, and TD3 [57, 58] learn a Q-function parameterized by  $\theta$  that estimates the expected total discounted reward given a state and an action. The Q function can be computed iteratively as  $Q^\pi(s_t, a_t) = \mathbb{E}_{s_{t+1} \sim P_T} [R(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi} (Q^\pi(s_{t+1}, a_{t+1}))]$ , such that it minimizes the Bellman error over all sampled transitions. The deterministic policy can then be updated through the deterministic policy gradient [57]. However, the standard formulation faces two major challenges; high sample complexity and brittle convergence. To this end, Soft Actor-Critic (SAC) [19] proposes to alter the original objective to include entropy maximization, i.e.,  $J(\pi) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho^\pi} [R(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))]$ , where  $\alpha$  is a temperature parameter that tunes the stochasticity of the optimal policy and  $\mathcal{H}$  represents the entropy of the policy. The Q-function and value function can be computed iteratively by:

$$Q^\pi(s_t, a_t) = R(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P_T} [V^\pi(s_{t+1})] \quad (1)$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi} [Q^\pi(s_t, a_t) - \alpha \log \pi(a_t|s_t)] \quad (2)$$

The last term adds a larger value to the original Q function if the probability of taking action  $a$  is small so that it explores actions that are less likely to occur. By adding entropy, it improves exploration, allowing for searching sub-optimal actions that may potentially lead to better rewards. The policy can then be trained to minimize the KL divergence  $J_\pi(\phi) = \mathbb{E}_{s_t \sim \rho^\pi} [\text{KL}(\pi_\phi(\cdot|s_t) \parallel \exp(Q_\theta^{\text{old}}(s_t, \cdot)) / Z_\theta^{\text{old}})]$  where  $Z^{\text{old}}$  is a partition function that normalizes the distribution.

## 4 Soft Hierarchical Reinforcement Learning

While HIRO is very effective, it searches the same areas redundantly (see Section 5), which contributes to data inefficiency, slowing down learning. In our algorithm, we propose to add entropy to the high-level policy to improve exploration. Furthermore, we also show how it can be best added to both levels. In Maximum Entropy RL, e.g., the Soft Actor Critic algorithm, the objective is to maximize the expectation of accumulated rewards and the policy entropy for a single agent. Similarly, in SHIRO, we also maximize entropy as follows

$$J(\pi^h) = \sum_{t \in T_c} \mathbb{E}_{(s_t, g_t) \sim \rho^{\pi^h}} [R^{abs}(s_t, g_t) + \alpha \mathcal{H}(\pi^h(\cdot|s_t, g_t))] \quad (3)$$

where we define an Abstracted Reward Function  $R^{abs}(s_t, g_t) = \sum_{t'=t}^{t+c} R(s_{t'}, \pi^l(a_{t'}|s_{t'}, g_t))$ . In the case of goal-conditioned HRL, the high-level agent provides a new sub-goal every  $c$  steps, where  $T_c$  is a set of time step  $t$  at every  $c$  steps, defined as  $T_c = \{t|t = c \cdot n, n = 0, 1, 2, \dots\}$ . By introducing such a function, the objective can become analogous to the standard reinforcement learning objective if we assume that the low-level policy is a part of the environment. Given this assumption, in order to transform the HRL problem into one that can be applied to standard RL, we introduce another function called the Abstracted Transition Function  $P_T^{abs}(s_{t+c}|s_t, g_t)$ ; i.e., a transition function from a state  $s_t$  to a state  $s_{t+c}$ , given a sub-goal  $g_t$ . With the abstraction, the probability of a trajectory  $\tau$  can be rewritten as follows, where the original probability for the single policy is given as a multiplication of probabilities of state and action sequences along the trajectory,  $p(\tau) = p(s_0) \prod_{t \in T_c} [\pi^h(g_t|s_t, g) P_T^{abs}(s_{t+c}|s_t, g_t)]$ . Given that the low-level policy gets updated, which induces a certain degree of stochasticity, we can view this Abstracted Transition Function as a representation of the stochastic environment. This is exactly the same problem as in the standard Reinforcement Learning objective for the high-level policy, as long as the Abstracted Transition

function does not change drastically. Formally, we can bound the change in the Abstracted Transition function by the following theorem:

**Theorem 1.** *Let  $\pi_\phi^l(a_t|s_t, g_t)$  be the low-level policy before the parameter update and  $\pi_{\phi'}^l(a_t|s_t, g_t)$  be the updated low-level policy after  $c$  steps. If two policies are **close**, i.e.,  $|\pi_{\phi'}^l(a_t|s_t, g_t) - \pi_\phi^l(a_t|s_t, g_t)| \leq \epsilon$  for all  $s_t$ , then the change in the abstracted transition functions is bounded by a small value  $|P_{T\pi_{\phi'}^l}^{abs}(s_{t+c}|s_t, g_t) - P_{T\pi_\phi^l}^{abs}(s_{t+c}|s_t, g_t)| \leq 2\epsilon c$ , where  $\epsilon \in [0, 1]$ .*

See the proof in Appendix A.

#### 4.1 Addition of Entropy to Each Policy

**High-Level Policy:** If the changes in the parameters are small, and if we can assume the low-level policy as a part of the environment, we can view the HRL problem in the same way as in the original RL problem, shown in Equation (3). Then, we can derive the exact same theorems as in SAC [19] for the high-level policy. For a deterministic high-level, we can apply the theorems of standard policy improvement and iteration. The addition of entropy to high-level is desirable for improved hierarchical exploration because it will generate more diverse sub-goals that facilitate learning of the low-level, leading to a better performing low-level agent to reach those broad sub-goals.

**Low-Level Policy:** Simultaneously (or by itself), we can add entropy to the low-level policy according to Theorem 1. To make the total variation divergence of policies between the parameter update small, we can use Pinsker’s inequality, i.e.,  $|\pi_{\phi'}^l(a_t|s_t, g_t) - \pi_\phi^l(a_t|s_t, g_t)| \leq \sqrt{\frac{1}{2}\text{KL}(\pi_{\phi'}^l(a_t|s_t, g_t)||\pi_\phi^l(a_t|s_t, g_t))}$ . This means that if the KL divergence between parameter updates is kept low, we can bound the changes in the Abstract Transition function and allow the low-level policy to be part of the environment. In this paper, we empirically show that the high-level and both-levels agents can perform better than others, and we will verify that the KL divergence is sufficiently small during training in Section 5. Moreover, the temperature parameter should not be too large to avoid the second term in Equation (2) to be large, which induces an incorrect estimation of the value function. This induces overly high values on actions that are unlikely to be taken, messing up the estimation of the low-level policy.

#### 4.2 Learning Temperature Parameter

Given that our proof allows to view the high-level agent as a single SAC agent, this motivates the use of learned entropy on the high-level. We apply the learned temperature method of [59] into the high-level and low-level agent. The objective for learning entropy temperature is given by  $\arg \min_{\alpha} \mathbb{E}_{a_t \sim \pi_t^*} [-\alpha \log \pi_t^*(a_t|s_t, g_t) - \alpha \mathcal{H}]$ , where  $\mathcal{H}$  is the minimum desired entropy, and  $\pi_t^*$  is the policy at timestep  $t$ . The expression is a measure of the difference between the current policy’s entropy and the minimum desired entropy. A lower value of  $\alpha$  will minimize and scale down this difference. With respect to the high-level, we assign a high temperature at the beginning of training to speed up learning, but allowing the temperature to decrease once the agent has learned. Furthermore, this is sufficient as the low-level agent will have been sufficiently trained on diverse data initially; lowering entropy does not pose any major concern. Similarly, we can also learn the temperature parameter  $\alpha$  of the low-level agent, in which we desire a low entropy that won’t increase significantly.

## 5 Experiments

We conducted two experiments to demonstrate the performance of our algorithms. The first experiment compares our best method to HIRO in benchmark environments. The second experiment is an ablation analysis of our algorithm. In this study, entropy is added to either or both of the levels, as well as with learning and not learning the temperature to analyze the behavior of the algorithm. These experiments demonstrate the power of the maximum entropy hierarchy with a simulated robot in MuJoCo.

Concretely, the agent is modeled as in Figure 1. The structures and the hyperparameters of the policies and Q-functions in the deterministic case are exactly the same as in [8] to compare against HIRO fairly. High-level and low-level agents are trained simultaneously, while the high-level agent is

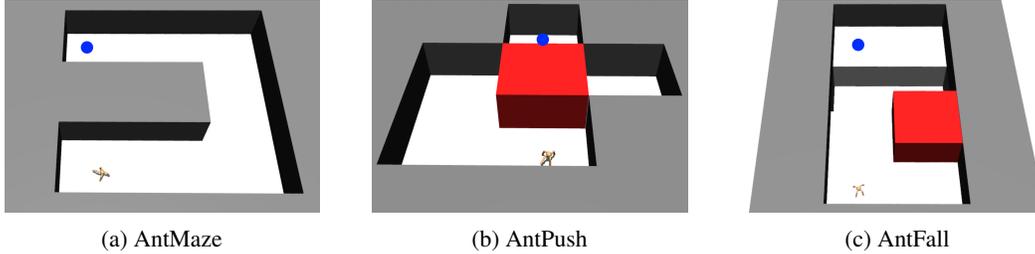


Figure 2: Three environments used in the experiments. In each environment, the ant starts at the initial location and aims to reach the goal colored in blue.

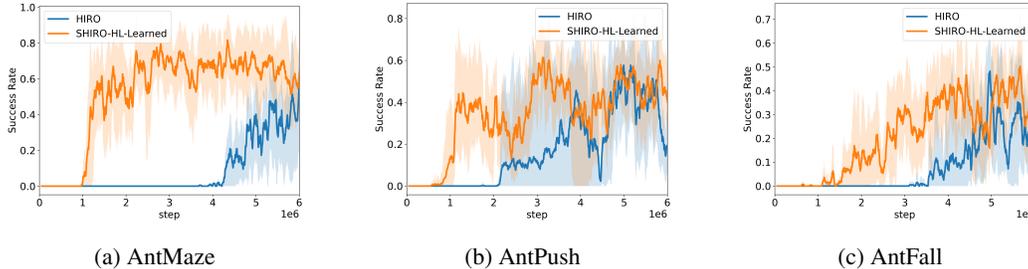


Figure 3: Success rate across different environments comparing SHIRO HL-Learned to HIRO. The darker lines represent the mean and the shaded areas represent the minimum and the maximum values across 3 different seeds.

updated every  $c = 10$  environment steps whereas the low-level agent is updated every environment step. The training of each agent is analogous to training a single agent. They can be trained independently using the samples from the replay buffer that are simultaneously relabeled by the off-policy correction method during the training. Most importantly, the entropy is defined and added to the original objective as in Equation (3). The entropy is added to either or both of the policies, and they are compared against HIRO, as in Figure 3. When we maximize entropy on either layer, we use the parameters according to Haarnoja’s ([19]) SAC. We call each of the variants, SHIRO HL (high-level entropy), SHIRO LL (low-level entropy), and BL (entropy at both levels) and used a temperature of 1.0 for high-level policy and 0.1 for the low-level policy across all environments. Other hyperparameters were carried over from SAC. We also used the same value as the initial temperatures for agents with temperature learning (SHIRO HL-Learned, SHIRO LL-Learned, SHIRO BL-Learned). All of the implementations is built on top of PFRL in PyTorch [60].

## 5.1 Comparative Study

In this study, we compared our *best* performing method, SHIRO-HL Learned, against HIRO [8] to demonstrate the effectiveness of entropy maximization formulation. HIRO is known to perform better than other HRL methods such as FeUdalNet[43] and Option-Critic[30] on the same tasks [8]. Furthermore, we also compared our methods to the vanilla SAC (no hierarchy).

In our research, we used the same environment that was used in the original paper of HIRO [8] to compare against their algorithm. We tested our algorithm on three of their environments: AntMaze, AntPush, and AntFall, as shown in Figure 2. All three environments require robotic locomotion and are standard benchmarks of HRL. As presented in Figure 3, in all environments, our method starts succeeding two to four times faster than the standard HIRO. Although the success rate increases steadily for AntMaze, it fluctuates for AntPush and AntFall, but gradually increases after the drop. This is thought of as a result of exploration for a better policy: Entropy maximization fosters exploration in new regions. The effect of broad exploration is shown in Figures 4a and 4b. The figure represents the scatter and contour plot of the final positions achieved by an agent in each environment. From the figure, we can observe that HIRO agent is failing to explore the areas around the goal but repetitively searches around the starting position before eventually learning locomotion to achieve the goal. In comparison to HIRO, qualitatively, our method is being able to reach goals more broadly.

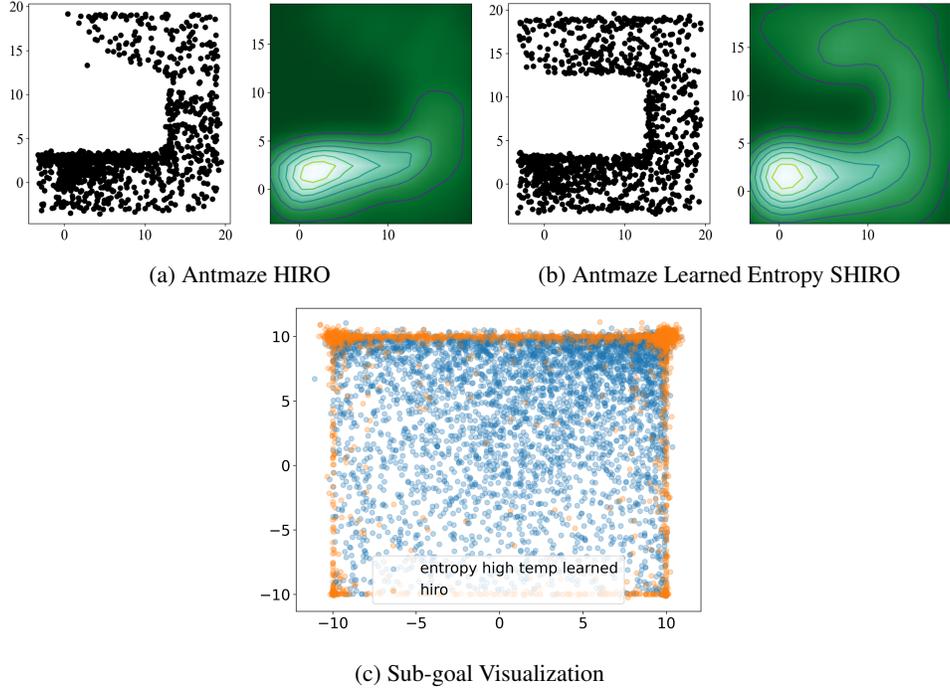


Figure 4: Top: Analysis of the final positions achieved by an agent in a given environment. The left plot is simply a scatterplot of the final positions of an agent on the environment specifies underneath. The second plot displays a contour map showing where the agent ends up on average across the entire experiment. We find that across all environments tested on, SHIRO HL-Learned explores the environment more thoroughly than HIRO. Bottom: X-Y Sub-goal visualization of the high-level policy.

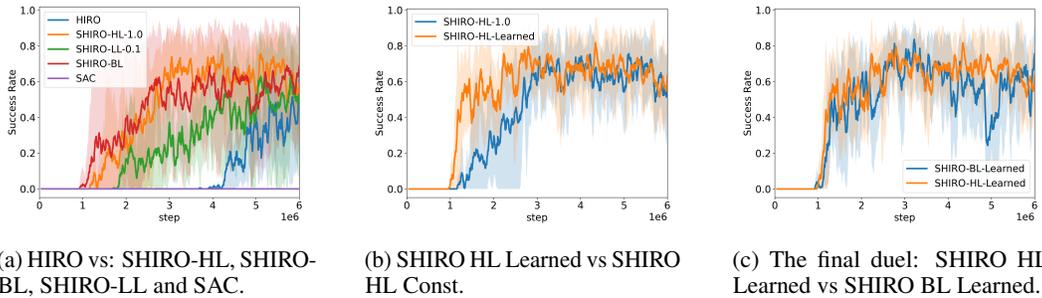


Figure 5: Ablative Study conducted in AntMaze. The center lines represent the mean and the shades represent the minimum and maximum range.

## 5.2 Ablative Study

In Figure 5, we display the results of our method SHIRO with different configurations that we tested. There are a few different ways to add entropy to a hierarchical agent, so we performed an ablative study to find out the best configuration.

**Adding Entropy to Each Layer.** We compare agents with entropy on different levels (SHIRO HL/LL/BL), as well as SAC. Each run result consists of three experiments using the same randomly selected seeds for each ablation. Figure 5a shows our results. We found that while all ablations perform better than the standard HIRO (besides SAC, which failed to complete the task), some are better than others. Adding entropy to the high-level controller yielded results that were consistently better across the board. Adding entropy to the low-level controller didn't perform as well, starting to succeed faster than regular HIRO, but took longer to reach the maximum success rate (Details in

Section 5.3). Finally, adding entropy to both levels performed at around the same level as adding it to just the high-level controller.

We find that the learned temperature parameter plays a role in the performance of our ablations, as in Figure 5b. If the temperature is high, the agent is incentivized to select actions (sub-goals in the case of the high-level agent) that provide more entropy. With respect to the high-level, while this is good when the agent is starting to learn, we want the agent to start acting deterministically as its success rate reaches its peak. We can observe this behavior in one of the figures in Appendix D, where the high-level temperature generally decreases with time according to the performance of the agent.

**High-Level vs. Both-Levels** Since both SHIRO HL-Learned and BL-Learned perform well, it is necessary to compare them in order to find the best option, shown in Figure 5c. We find that SHIRO HL-Learned performs marginally better than SHIRO BL-Learned. Thus, we consider SHIRO HL-Learned to be our best method produced from the ablative study, as well as the fact that SHIRO-HL-Learned does not require any temperature parameter tuning on the low-level. However, we suggest that both methods are viable, and one or the other can be chosen according to preference.

### 5.3 Analysis

In our experiments, the KL-divergence in the low-level policy is sufficiently low in our ablations, which follows Theorem 1 that can be applied to HIRO with and without high-level entropy. In general, the KL divergence is below 1.0 (a plot can be found in Appendix D). However, even though SHIRO with a low-level temperature of 1.0 had the lowest KL divergence, it fails to learn because the high temperature parameter induces Q-value overestimation. As the temperature increases, the second term in Equation (2) becomes large that overestimates the actions that are unlikely to be taken. It also creates a too stochastic environment for the high-level agent to solve the main task.

We also observe that the addition of entropy to the high-level agent significantly improves the learning of the *low-level* agent, in turn, improving the high-level agent by making the environment more “helpful” to solve the main task, as well as being able to broadly explore the environment. The more diverse sub-goals provided to the low-level agent increase the loss (see Appendix D) of that agent’s Q-functions because the Q-functions must initially learn from more diverse data (this difference in diversity can be seen in Figure 4c). This allows for larger updates (from gradient descent) that enable the Q-function to learn from underrepresented goals. With a more robust Q-function that has received more meaningful learning signals, we also observe an improvement in the low-level policy loss (details in Appendix D), where the policy directly optimizes from the Q-function. In turn, the low-level agent learns locomotion quickly.

Another observation from our study is that high-level entropy maximization improves slightly “worse” low-level agents. When the high-level agent is kept deterministic, selecting a low-level policy as SAC (which has empirically shown to outperform TD3) demonstrates better performance than a TD3 low-level policy. However, when the high-level policy maximizes entropy, the performance of different low-level agents significantly improves but are similar to one another, suggesting that the diversity of data significantly improves worse low-level agents. Graphs of these observations are located in Appendix D.

## 6 Conclusion

In this work, we proposed an effective solution to the problem of redundant sub-goals for robotic learning in goal-conditioned hierarchical reinforcement learning. We achieved this through maximum entropy goal-conditioned HRL, in which we treated it as a single agent when the KL divergence is small between the low-level policy updates. Our empirical findings motivate the learning of entropy temperature, namely on the high-level. We performed an ablative study to understand the effects of entropy on both levels, where high-level learned entropy had significant impact on the efficiency. Concretely, the diversity of sub-goals provided to the low-level agent improved the redundancy and reachability of goals, allowing for a high-level agent to learn even faster. Regarding future work, we are interested in applying our findings on a real robot platform, as well as exploration of the importance of sub-goal representations.

## References

- [1] S. Gu, E. Holly, T. Lillicrap, and S. Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 3389–3396. IEEE, 2017.
- [2] Y. Yang, K. Caluwaerts, A. Iscen, T. Zhang, J. Tan, and V. Sindhwani. Data efficient reinforcement learning for legged robots. In *Conference on Robot Learning*, pages 1–10. PMLR, 2020.
- [3] D. Jain, A. Iscen, and K. Caluwaerts. Hierarchical reinforcement learning for quadruped locomotion. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7551–7557. IEEE, 2019.
- [4] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International conference on machine learning*, pages 1329–1338. PMLR, 2016.
- [5] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [6] X. B. Peng, M. Chang, G. Zhang, P. Abbeel, and S. Levine. Mcp: learning composable hierarchical control with multiplicative compositional policies. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 3686–3697, 2019.
- [7] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [8] O. Nachum, S. S. Gu, H. Lee, and S. Levine. Data-efficient hierarchical reinforcement learning. *Advances in Neural Information Processing Systems*, 31:3303–3313, 2018.
- [9] A. Levy, G. Konidaris, R. Platt, and K. Saenko. Learning multi-level hierarchies with hindsight. In *International Conference on Learning Representations*, 2018.
- [10] Y. Jiang, S. S. Gu, K. P. Murphy, and C. Finn. Language as an abstraction for hierarchical deep reinforcement learning. *Advances in Neural Information Processing Systems*, 32:9419–9431, 2019.
- [11] O. Nachum, S. Gu, H. Lee, and S. Levine. Near-optimal representation learning for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2018.
- [12] P. Dayan and G. E. Hinton. Feudal reinforcement learning. In S. Hanson, J. Cowan, and C. Giles, editors, *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann, 1993.
- [13] R. Parr and S. Russell. Reinforcement learning with hierarchies of machines. *Advances in neural information processing systems*, pages 1043–1049, 1998.
- [14] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- [15] T. G. Dietterich. Hierarchical reinforcement learning with the maxq value function decomposition. *Journal of artificial intelligence research*, 13:227–303, 2000.
- [16] O. Nachum, H. Tang, X. Lu, S. Gu, H. Lee, and S. Levine. Why does hierarchy (sometimes) work so well in reinforcement learning? *arXiv preprint arXiv:1909.10618*, 2019.
- [17] T. Zhang, S. Guo, T. Tan, X. Hu, and F. Chen. Generating adjacency-constrained subgoals in hierarchical reinforcement learning. *Advances in Neural Information Processing Systems*, 33, 2020.
- [18] S. Li, L. Zheng, J. Wang, and C. Zhang. Learning subgoal representations with slow dynamics. In *International Conference on Learning Representations*, 2021.

- [19] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- [20] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- [21] M. Toussaint. Robot trajectory optimization using approximate inference. In *Proceedings of the 26th annual international conference on machine learning*, pages 1049–1056, 2009.
- [22] K. Rawlik, M. Toussaint, and S. Vijayakumar. On stochastic optimal control and reinforcement learning by approximate inference. *Proceedings of Robotics: Science and Systems VIII*, 2012.
- [23] R. Fox, A. Pakman, and N. Tishby. Taming the noise in reinforcement learning via soft updates. In *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, pages 202–211, 2016.
- [24] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pages 1352–1361. PMLR, 2017.
- [25] S. Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- [26] M. Stolle and D. Precup. Learning options in reinforcement learning. In *International Symposium on abstraction, reformulation, and approximation*, pages 212–223. Springer, 2002.
- [27] S. Mannor, I. Menache, A. Hoze, and U. Klein. Dynamic abstraction in reinforcement learning via clustering. In *Proceedings of the twenty-first international conference on Machine learning*, page 71, 2004.
- [28] S. Singh, A. G. Barto, and N. Chentanez. Intrinsically motivated reinforcement learning. Technical report, MASSACHUSETTS UNIV AMHERST DEPT OF COMPUTER SCIENCE, 2005.
- [29] D. Precup. *Temporal abstraction in reinforcement learning*. University of Massachusetts Amherst, 2000.
- [30] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17*, page 1726–1734. AAAI Press, 2017.
- [31] J. Harb, P.-L. Bacon, M. Klissarov, and D. Precup. When waiting is not an option: Learning options with a deliberation cost. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [32] A. Vezhnevets, V. Mnih, S. Osindero, A. Graves, O. Vinyals, J. Agapiou, et al. Strategic attentive writer for learning macro-actions. *Advances in Neural Information Processing Systems*, 29: 3486–3494, 2016.
- [33] K. Frans, J. Ho, X. Chen, P. Abbeel, and J. Schulman. Meta learning shared hierarchies. In *International Conference on Learning Representations*, 2018.
- [34] O. Sigaud and F. Stulp. Policy search in continuous action domains: an overview. *Neural Networks*, 113:28–40, 2019.
- [35] N. Heess, G. Wayne, Y. Tassa, T. Lillicrap, M. Riedmiller, and D. Silver. Learning and transfer of modulated locomotor controllers. *arXiv preprint arXiv:1610.05182*, 2016.
- [36] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29:3675–3683, 2016.
- [37] C. Tessler, S. Givony, T. Zahavy, D. J. Mankowitz, and S. Mannor. A deep hierarchical approach to lifelong learning in minecraft. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17*, page 1553–1561. AAAI Press, 2017.

- [38] C. Florensa, Y. Duan, and P. Abbeel. Stochastic neural networks for hierarchical reinforcement learning. In *ICLR (Poster)*, 2017.
- [39] G. D. Konidaris and A. G. Barto. Building portable options: Skill transfer in reinforcement learning. In *IJCAI*, volume 7, pages 895–900, 2007.
- [40] C. Daniel, G. Neumann, and J. Peters. Hierarchical relative entropy policy search. In *Artificial Intelligence and Statistics*, pages 273–281. PMLR, 2012.
- [41] T. Schaul, D. Horgan, K. Gregor, and D. Silver. Universal value function approximators. In *International conference on machine learning*, pages 1312–1320. PMLR, 2015.
- [42] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba. Hindsight experience replay. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 5055–5065, 2017.
- [43] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International Conference on Machine Learning*, pages 3540–3549. PMLR, 2017.
- [44] O. Nachum, M. Ahn, H. Ponte, S. S. Gu, and V. Kumar. Multi-agent manipulation via locomotion using hierarchical sim2real. In *Conference on Robot Learning*, pages 110–121. PMLR, 2020.
- [45] S. Mahadevan and M. Maggioni. Proto-value functions: A laplacian framework for learning representation and control in markov decision processes. *Journal of Machine Learning Research*, 8(10), 2007.
- [46] R. S. Sutton, J. Modayil, M. Delp, T. Degris, P. M. Pilarski, A. White, and D. Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pages 761–768, 2011.
- [47] V. Pong, S. Gu, M. Dalal, and S. Levine. Temporal difference models: Model-free deep rl for model-based control. In *International Conference on Learning Representations*, 2018.
- [48] A. Li, C. Florensa, I. Clavera, and P. Abbeel. Sub-policy adaptation for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2019.
- [49] D. Hejna, L. Pinto, and P. Abbeel. Hierarchically decoupled imitation for morphological transfer. In *International Conference on Machine Learning*, pages 4159–4171. PMLR, 2020.
- [50] J. Zhang, H. Yu, and W. Xu. Hierarchical reinforcement learning by discovering intrinsic options. In *International Conference on Learning Representations*, 2020.
- [51] S. Li, R. Wang, M. Tang, and C. Zhang. Hierarchical reinforcement learning with advantage-based auxiliary rewards. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 1409–1419, 2019.
- [52] R. Zhao, X. Sun, and V. Tresp. Maximum entropy-regularized multi-goal reinforcement learning. In *International Conference on Machine Learning*, pages 7553–7562. PMLR, 2019.
- [53] T. Haarnoja, K. Hartikainen, P. Abbeel, and S. Levine. Latent space policies for hierarchical reinforcement learning. In *International Conference on Machine Learning*, pages 1851–1860. PMLR, 2018.
- [54] A. Azarafrooz and J. Brock. Hierarchical soft actor-critic: Adversarial exploration via mutual information optimization. *arXiv preprint arXiv:1906.07122*, 2019.
- [55] H. Tang, A. Wang, F. Xue, J. Yang, and Y. Cao. A novel hierarchical soft actor-critic algorithm for multi-logistics robots task allocation. *IEEE Access*, 9:42568–42582, 2021.
- [56] B. Beyret, A. Shafti, and A. A. Faisal. Dot-to-dot: Explainable hierarchical reinforcement learning for robotic manipulation. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5014–5019. IEEE, 2019.

- [57] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In *ICLR (Poster)*, 2016.
- [58] S. Fujimoto, H. Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596. PMLR, 2018.
- [59] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- [60] Y. Fujita, P. Nagarajan, T. Kataoka, and T. Ishikawa. Chainerrl: A deep reinforcement learning library. *Journal of Machine Learning Research*, 22(77):1–14, 2021. URL <http://jmlr.org/papers/v22/20-376.html>.
- [61] A. Abdolmaleki, J. T. Springenberg, Y. Tassa, R. Munos, N. Heess, and M. Riedmiller. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations*, 2018.

## A Appendix A: Proofs and Additional Derivations

### A.1 Proof of Theorem 1

*Proof.* In this proof, we show that the change in the Abstract Transition function  $P_{T\pi_\phi}^{abs}(s_{t+c}|s_t, g_t)$  can be bounded if the policy is updated such that the policies before and after the parameter update after time step  $c$  are *close*. In the deterministic case, we say that the policy  $\pi_{\phi'}$  after the update is *close* to the policy  $\pi_\phi$  before the update if taking different action from the old policy is bounded by a small probability value  $\pi_{\phi'}^l(a_t \neq \pi_\phi^l(s_t, g_t)|s_t, g_t) \leq \epsilon$ . Let the probability of arriving at state  $s_{t+c}$  starting from state  $s_t$  given sub-goal  $g_t$  be

$$p_{\phi'}(s_{t+c}|s_t, g_t) = (1 - \epsilon)^c p_\phi(s_{t+c}|s_t, g_t) + (1 - (1 - \epsilon)^c) p_{\text{mistakes}}(s_{t+c}|s_t, g_t) \quad (4)$$

The  $p_{\phi'}(s_{t+c}|s_t, g)$  is defined as an addition of the probability of not making mistake for  $c$  consecutive steps and the probability of making at least one mistake, where  $p_{\text{mistakes}}$  is some distribution that makes mistakes. Then the difference in the probability, defined as total variation divergence between  $p_{\phi'}$  and  $p_\phi$ , is bounded by  $2\epsilon c$ .

$$|p_{\phi'}(s_{t+c}|s_t, g_t) - p_\phi(s_{t+c}|s_t, g_t)| = (1 - (1 - \epsilon)^c) |p_{\text{mistakes}}(s_{t+c}|s_t, g_t) - p_\phi(s_{t+c}|s_t, g_t)| \quad (5)$$

$$\leq 2(1 - (1 - \epsilon)^c) \quad (6)$$

$$\leq 2\epsilon c \quad (7)$$

Thus, as long as the probability of a policy making a mistake is small, then the probability of transitioning to  $s_{t+c}$  is almost the same as before. Notice that  $p_\phi(s_{t+c}|s_t, g_t)$  is exactly equal to  $P_{T\pi_\phi}^{abs}(s_{t+c}|s_t, g_t)$ . Thus, we can say that the change in the Abstracted Transition function does not change much if policy stays similar. In fact, we can similarly derive a bound for the stochastic case as in [25]. As long as the total variation divergence between the two policies is bounded by  $\epsilon$ , i.e.,  $|\pi_{\phi'}^l(a_t|s_t, g_t) - \pi_\phi^l(a_t|s_t, g_t)| \leq \epsilon$  for all  $s_t$ , then the policies are said to be *close*.  $\square$

### A.2 KL Divergence of the Deterministic Policy with Gaussian Noise

In our work, we computed the KL divergence between low-level policies to compare the results against other methods and to verify our theoretical findings. For deterministic policies, especially policies with Gaussian exploration noise, we view them as stochastic and derive the KL divergence as follows. Let  $\mu : \mathcal{S} \times \mathcal{G} \rightarrow \mathcal{A}$  be the deterministic head of the policy and  $\Sigma$  be the exploration noise.

$$\pi(a_t|s_t, g_t) = \mathcal{N}(\mu(s_t, g_t), \Sigma), \quad (8)$$

Thus, under this perspective, the KL divergence between a old policy  $\pi_\phi$  and new  $\pi_{\phi'}$  after a parameter update can be expressed as:

$$D_{KL}(\pi_\phi || \pi_{\phi'}) = \frac{1}{2} \left[ \log \frac{|\Sigma_{\phi'}|}{|\Sigma_\phi|} - k + (\mu_{\phi'} - \mu_\phi)^T \Sigma_{\phi'}^{-1} (\mu_{\phi'} - \mu_\phi) + \text{tr}\{\Sigma_{\phi'}^{-1} \Sigma_\phi\} \right], \quad (9)$$

where  $\Sigma_\phi$  is the covariance matrix of  $\pi_\phi$  before the parameter update,  $\Sigma_{\phi'}$  is the covariance matrix after the parameter update,  $k$  is the action dimension, and  $\mu_\phi$  and  $\mu_{\phi'}$  are the means of the policies before and after the parameter update, respectively. In our implementation, we set the covariance of  $\Sigma = \sigma I$  where  $\sigma$  is a constant positive value and  $I$  is an identity matrix. Then, we can simply the notation as,

$$D_{KL}(\pi_\phi || \pi_{\phi'}) = \frac{1}{2} (\mu_{\phi'} - \mu_\phi)^T \Sigma_{\phi'}^{-1} (\mu_{\phi'} - \mu_\phi), \quad (10)$$

## B Appendix B: Environment

In all environments, the goal of the agent is to reach the desired goal position, regardless of obstacles. The agent constantly receives a scalar reward that is defined as a negative distance from its current position to the desired goal, i.e.,  $R(s_t, a_t) = -\sqrt{\sum_i |g_i - s_i|^2}$  where  $i \in \{0, 1\}$  for 2D and  $\{0, 1, 2\}$  for 3D goals. As the simulated robot gets closer to the goal, the reward increases. AntMaze is a maze in a shape of  $\supset$  where an agent starts from the bottom left at location (0, 0) and aims to reach the top left goal at location of (0, 16). This experiment tests whether an algorithm can avoid one of the local optima to search for a better policy. Next, AntPush is a maze where a large box is placed in the center of the environment to block the agent from reaching the final goal at location (0, 19). The agent’s goal is to move the block so that its path to the goal is cleared out. Similarly, in AntFall, the agent has to push a block into a ditch to cross over the ditch to get to its other end to reach the final goal at location (0, 27, 4.5). This environment extends the navigation task to three dimensions. In both AntPush and AntFall, the greedy algorithm will fail to push the block first, but rather takes a greedy action and become stuck in a spot with decent rewards but not the higher possible reward. Across all experiments, a four-legged Ant is used as a robotic platform. Each leg has two limbs and each limb has a joint. The ant controls all of its joints to walk itself to its desired position. The control input is consisted of these eight inputs. It then observes 31 dimensional states that consists of ant location (3), ant orientation (4), joint angles (8), and all of their velocities, plus time step  $t$ . The limit of each state and control input are detailed as follows.

The limit of the control input space  $\pm[30, 30, 30, 30, 30, 30, 30, 30]$

The limit of the sub-goal space  $\pm[10, 10, 0.5, 1, 1, 1, 1, 0.5, 0.3, 0.5, 0.3, 0.5, 0.3, 0.5, 0.3]$

## C Appendix C: Implementation Details

We referenced the implementation of HIRO and used the same architecture and hyperparameters. We adopted the basic TD3 architecture for the deterministic policy in both levels. As reported in HIRO, we also modified the layer size of (400, 300) to (300, 300). In both layers, we mostly used the same hyperparameters: discount factor  $\gamma = 0.99$ , actor learning rate 0.0001, critic learning rate 0.001, soft update targets  $\tau = 0.005$ , replay buffer of size 200,000, no gradient clipping, and Gaussian noise with  $\sigma = 1.0$  for exploration. For the reward scaling, 0.1 was used for the high-level and 1.0 for the low-level. The training step and target update interval were set to 1 environment step for low-level and 10 environment steps for high-level. Lastly, the time step  $c$  is set to 10 environment steps.

Soft Constraint:

$$\max_{\phi} J(\pi^l) = \max_{\phi} \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho^{\pi^l}} \left[ r(s_t, a_t) - \alpha_{KL} \cdot \text{KL}(\pi_{\phi'}^1(\cdot | s_t) \| \pi_{\phi}^1(\cdot | s_t)) \right] \quad (11)$$

For the soft constraint, we can add a KL divergence loss to the original policy loss when updating the current policy, i.e.,

$$L'_{\phi}(\phi) = L_{\phi}(\phi) + \alpha_{KL} \cdot \mathbb{E}_{s_t \sim \rho^{\pi^l}} \left[ \text{KL}(\pi_{\phi'}^1(\cdot | s_t) \| \pi_{\phi}^1(\cdot | s_t)) \right] \quad (12)$$

Hard Constraint: E-Step

$$\max_{\phi} J(\pi^l) = \max_{\phi} \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho^{\pi^l}} \left[ r(s_t, a_t) \right] \quad (13)$$

$$\text{s.t.} \quad \mathbb{E}_{s_t \sim \rho^{\pi^l}} \left[ \text{KL}(\pi_{\phi'}^1(\cdot | s_t) \| \pi_{\phi}^1(\cdot | s_t)) \right] \leq \epsilon \quad (14)$$

In order to strictly constrain the KL divergence during the policy update, we need to derive an alternative method for optimizing the policy. We referred to Maximum a Posteriori Policy Optimisation (MPO) algorithm [61] that optimizes policy with the family of Expectation Maximization (EM) algorithm.

---

**Algorithm 1** SHIRO Psuedo-Code

---

Initialize Critic Networks  $Q_{\theta^{h_1}}, Q_{\theta^{h_2}}, Q_{\theta^{l_1}}, Q_{\theta^{l_2}}$  with random parameters  $\theta^{h_1}, \theta^{h_2}, \theta^{l_1}, \theta^{l_2}$   
Initialize Target Networks with parameters  $\theta^{h_1'} \leftarrow \theta^{h_1}, \theta^{h_2'} \leftarrow \theta^{h_2}, \theta^{l_1'} \leftarrow \theta^{l_1}, \theta^{l_2'} \leftarrow \theta^{l_2}$   
Initialize Actor Networks  $\pi^h$  and  $\pi^l$  with random parameters  $\phi^h, \phi^l$   
Initialize Replay Buffers  $\beta^h$  and  $\beta^l$   
Get initial environment state  $s_0$  and goal  $g$   
Temperature parameters  $\alpha^h, \alpha^l$   
**for** each environment step  $t$  **do**  
  Select action  $a_t \sim \pi^l(s_t, g_t)$  and get new observation and reward  $s_{t+1}, r_t$   
  **if** step % c == 0 **then**  
    Select sub-goal  $g_{t+1} \sim \pi^h(s_t, g_t)$   
  **else**  
     $g_{t+1} \leftarrow s_t + g_t - s_{t+1}$   
  **end if**  
  Store transition tuple  $(s_t, g_t, a_t, r_t, s_{t+1}, g_{t+1}, done_t)$  in  $\beta^l$   
  **if** step % low\_train\_freq == 0 **then**  
    sample mini-batch of N transitions from  $\beta^l$   
    **if**  $\alpha^l$  **then**  
      Entropy  $e \sim \pi_{\phi}^l(s_{t+1}, g_{t+1})$   
    **else**  
       $e = 0$   
    **end if**  
     $Q_{target}^i \leftarrow r_t + \gamma \cdot (1 - done) \cdot (Q_{\theta^{l_i}}^i(s_{t+1}, g_{t+1})) - \alpha e$   
     $Q_{predict}^i \leftarrow Q_{\theta^{l_i}}(s_t, g_t)$   
     $\theta^{l_i} \leftarrow \theta^{l_i} - \eta_{\theta} \nabla \mathcal{L}(Q_{target}^i, Q_{predict}^i)$   
    **if** step % low\_update\_delay == 0 **then**  
       $\phi^l \leftarrow \phi^l - \eta_{\phi} \nabla \mathbb{E}[Q_{\theta^{l_1}}(s_t, g_t) - \alpha e]$   
       $\theta^{l_i'} \leftarrow \tau \theta^{l_i} + (1 - \tau) \theta^{l_i'}$   
       $\phi^{l'} \leftarrow \tau \phi^l + (1 - \tau) \phi^{l'}$   
      **if**  $\alpha^l$  **then**  
         $\alpha^l \leftarrow \alpha^l - \eta_{\alpha} \nabla \mathbb{E}[\alpha \pi_{\phi}^l]$   
      **end if**  
    **end if**  
  **end if**  
  **if** step % high\_train\_freq == 0 **then**  
    Store transition tuple  $(s_{t-N:t}, g_{t-N:t}, a_{t-N:t}, R_{t-N:t}, s_{t+1}, done_{t-N:t})$  in  $\beta^h$   
    sample mini-batch of N transitions from  $\beta^h$   
    **if**  $\alpha^h$  **then**  
      Entropy  $e \sim \pi_{\phi}^h(s_{t+1}, g_{t+1})$   
    **else**  
       $e = 0$   
    **end if**  
     $Q_{target}^i \leftarrow R_t + \gamma \cdot (1 - done) \cdot (Q_{\theta^{h_i}}^i(s_{t+1}, g_{t+1})) - \alpha e$   
     $Q_{predict}^i \leftarrow Q_{\theta^{h_i}}(s_t, g_t)$   
     $\theta^{h_i} \leftarrow \theta^{h_i} - \eta_{\theta} \nabla \mathcal{L}(Q_{target}^i, Q_{predict}^i)$   
    **if** step % high\_update\_delay **then**  
       $\phi^h \leftarrow \phi^h - \eta_{\phi} \nabla \mathbb{E}[Q_{\theta^{h_1}}(s_t, g_t) - \alpha e]$   
       $\theta^{h_i'} \leftarrow \tau \theta^{h_i} + (1 - \tau) \theta^{h_i'}$   
       $\phi^{h'} \leftarrow \tau \phi^h + (1 - \tau) \phi^{h'}$   
      **if**  $\alpha^h$  **then**  
         $\alpha^h \leftarrow \alpha^h - \eta_{\alpha} \nabla \mathbb{E}[\alpha \pi_{\phi}^h]$   
      **end if**  
    **end if**  
  **end if**  
**end for**

---

## D Appendix D: Additional Analysis and Figures

We included additional figures for the analysis conducted in Section 5.3.

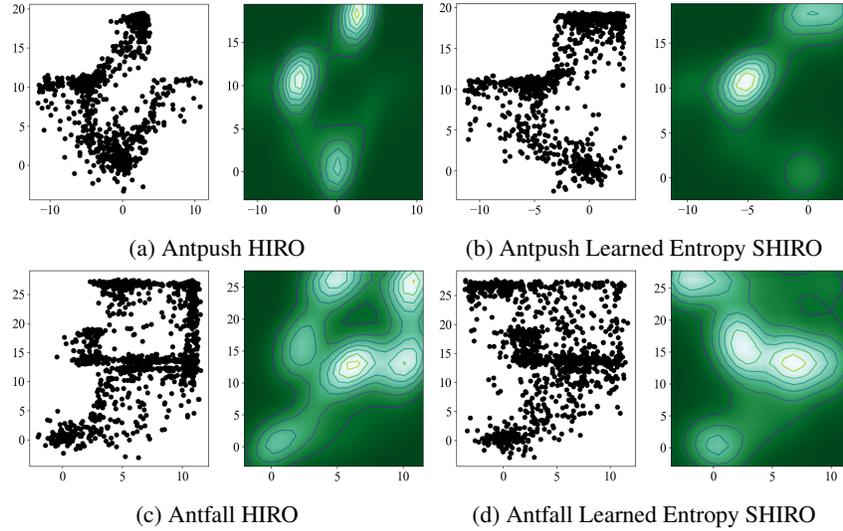


Figure 6: Analysis of the final positions achieved by an agent in a given environment. In each subfigure, the left plot is simply a scatterplot of the final positions of an agent on the environment specifies underneath. The second plot takes the final position data and displays a contour map, showing where the agent ends up on average across the entire experiment. We find that across all environments tested on, SHIRO HL-Learned explores the environment more thoroughly than HIRO.

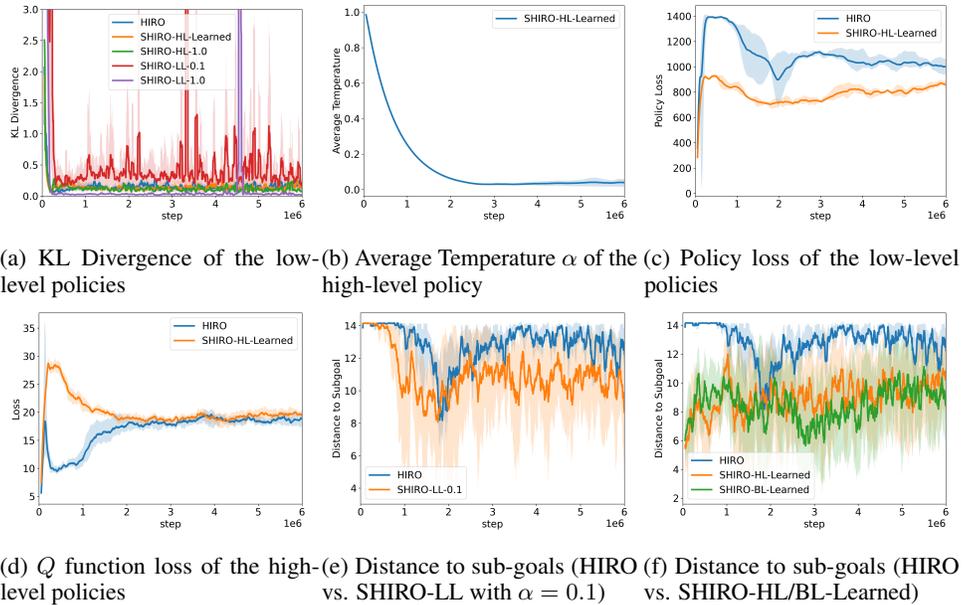


Figure 7: Additional figures for the analysis of the entropy effects.