

# Robustly Learning Composable Options in Deep Reinforcement Learning

Akhil Bagaria\*, Jason Senthil\*, Matthew Slivinski and George Konidaris

Department of Computer Science, Brown University

{akhil\_bagaria, jason\_senthil, matthew\_slivinski}@brown.edu, gdk@cs.brown.edu

## Abstract

Hierarchical reinforcement learning (HRL) is only effective for long-horizon problems when high-level skills can be reliably sequentially executed. Unfortunately, learning reliably composable skills is difficult, because all the components of every skill are constantly changing during learning. We propose three methods for improving the composability of learned skills: representing skill initiation regions using a combination of pessimistic and optimistic classifiers; learning re-targetable policies that are robust to non-stationary subgoal regions; and learning robust option policies using model-based RL. We test these improvements on four sparse-reward maze navigation tasks involving a simulated quadrupedal robot. Each method successively improves the robustness of a baseline skill discovery method, substantially outperforming state-of-the-art flat and hierarchical methods.<sup>1</sup>

## 1 Introduction

Temporal abstraction is a promising approach to scaling reinforcement learning (RL) algorithms [Botvinick *et al.*, 2009] because it is capable of addressing some of the biggest challenges in RL: structured exploration [Jinnai *et al.*, 2019], transfer [Konidaris and Barto, 2007] and long-term credit assignment [Dietterich, 2000]. Hierarchical methods are particularly attractive for long-horizon, sparse reward tasks, where flat (non-hierarchical) RL algorithms often struggle.

Hierarchical approaches succeed by allowing the agent to sequentially execute high-level actions. This intuition has led to several skill-discovery algorithms that explicitly satisfy the composability objective: executing one skill takes the agent to a state where it can execute another. Such algorithms are widespread in the literature; from control-theory [Lyapunov, 1992; Burrige *et al.*, 1999; Tedrake, 2009], to robotics [Lozano-Perez *et al.*, 1984; Kaelbling and Lozano-Pérez, 2017], to the online [Randløv *et al.*, 2000; Konidaris

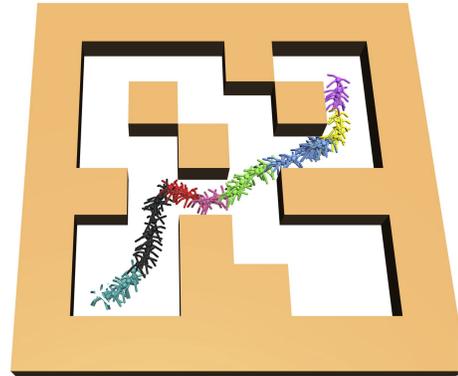


Figure 1: A learned solution to a long-horizon, sparse reward task that requires executing approximately 665 low-level steps (different colors denote discovered skills). Our proposed methods result in discovered skills that can be reliably sequentially executed, enabling skill discovery methods to more effectively solve such tasks.

and Barto, 2009; Shoeleh and Asadpour, 2017; Bagaria and Konidaris, 2020] and batch RL [Singh *et al.*, 2020] settings.

Even when skills are not explicitly constructed to be sequentially executable, they must eventually be sequenced to solve goal-directed tasks [Sharma *et al.*, 2020; Frans *et al.*, 2018]. Additionally, skills that can be reliably sequenced can support abstract, high-level planning [Konidaris *et al.*, 2018].

The core difficulty that arises when discovering composable skills is that of *non-stationary subgoals*. Skill-discovery algorithms that explicitly optimize for composability must learn the regions (called *initiation sets* [Sutton *et al.*, 1999]) from which each skill can be successfully executed. To construct sequentially executable skills, the initiation set of one skill becomes the subgoal of a new skill [Konidaris and Barto, 2009]. These initiation sets are constantly changing as the agent learns, causing the target (and therefore reward function) of successive skills to also change over time, destabilizing composability and complicating learning.

We propose three methods to combat this problem. First, we propose to stabilize learned initiation sets using a dual-classifier approach. An optimistic classifier determines when the agent can execute the skill, which encourages exploration. Meanwhile, a pessimistic classifier is used as a subgoal target, which is likely to grow but unlikely to shrink, leading to stable chains. Next, to create skill policies that are robust to

\*Equal contribution

<sup>1</sup>Video, code and appendix can be found at <https://sites.google.com/brown.edu/robustly-composing-options>

non-stationary subgoals, we propose to use goal-conditioned policies [Schaul *et al.*, 2015] for each skill. Such policies are robust to subgoal changes because they can always be re-targeted towards a state in the new subgoal region. Finally, while the majority of skill-discovery work has been model-free, we show that using model-based RL to learn skill policies leads to a substantially more robust skills.

To evaluate our proposed changes, we consider four sparse-reward continuous control problems in MuJoCo [Todorov *et al.*, 2012]. Compared to a flat model-free solution [Fujimoto *et al.*, 2018; Andrychowicz *et al.*, 2017], our agent achieves at least two orders of magnitude better sample efficiency. While a flat model-based solution [Nagabandi *et al.*, 2018] is unable to solve any of the problems considered in this paper, our model-based hierarchy solves them with relative ease. Finally, our agent achieves at least a  $5\times$  improvement in sample-efficiency over a state-of-the-art skill-discovery algorithm [Bagaria and Konidaris, 2020].

## 2 Background and Related Work

We consider decision making problems modelled as episodic, goal-oriented MDPs  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma, \mathcal{G})$  [Sutton and Barto, 2018], where  $\mathcal{G}$  refers to a set of goals that the agent could be asked to reach, whereupon execution terminates with a non-negative reward. Like most goal-conditioned RL algorithms [Schaul *et al.*, 2015; Andrychowicz *et al.*, 2017], we assume access to a reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{G} \rightarrow \mathbb{R}$ .

**Model-Free RL:** Model-free methods often use Q-learning [Watkins and Dayan, 1992] to learn  $Q^\pi(s_t, a_t)$ , estimating the expected sum of discounted future rewards conditioned on taking action  $a_t$  from  $s_t$ , thereafter following policy  $\pi$ . The policy  $\pi(s_t) = \arg \max_{a \in \mathcal{A}} Q(s_t, a)$  chooses an action that greedily maximizes the Q-function at a given state. The Q-function is often represented by a neural network  $\phi$  [Mnih *et al.*, 2015] and the policy by another neural network  $\psi$  [Lillicrap *et al.*, 2015]. Many algorithms exist for learning  $\phi$  and  $\psi$  from interactions with  $\mathcal{M}$  [Duan *et al.*, 2016].

**Hindsight Experience Replay (HER):** When the class of policies we wish to learn in  $\mathcal{M}$  is restricted to goal-reaching policies [Kaelbling, 1993], HER can improve sample efficiency as follows. Given a trajectory  $\tau_1 = (s_1, a_1, \mathcal{R}(s_2, g), s_2), \dots, (s_n, a_n, \mathcal{R}(s_{n+1}, g), s_{n+1})$  executed while trying to reach goal  $g$ , HER [Andrychowicz *et al.*, 2017] replays  $\tau_1$  [Lin, 1993] assuming that the agent was trying to reach goal  $s_{n+1}$ , by augmenting its set of experiences with another trajectory  $\tau_2 = (s_1, a_1, \mathcal{R}(s_2, s_{n+1}), s_2), \dots, (s_n, a_n, \mathcal{R}(s_{n+1}, s_{n+1}), s_{n+1})$ . Replaying trajectory  $\tau_2$  trains the agent to reach state  $s_{n+1}$ , even when the reward function  $\mathcal{R}$  is sparse.

**Model-Based RL:** While Q-learning directly learns the value function for control, model-based methods first learn an approximate model of the system dynamics. This model  $f_\xi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  is often represented using another neural network  $\xi$  [Nagabandi *et al.*, 2018]. Model-based methods can then select actions by solving:

$$\pi(s|g) = \arg \max_{a \in \mathcal{U}^{|\mathcal{A}|}(-1,1)} \sum_{t=1}^H \gamma^{t-1} R(s_t, g), \quad (1)$$

such that  $s_{t+1} = f_\xi(s_t, a)$ . In Model Predictive Control (MPC) [Garcia *et al.*, 1989], the agent approximately solves this optimization problem at every time-step, executing only the first action from the resulting action sequence. By re-planning at every time-step, MPC is robust to the prediction errors that compound over time with a learned model [Asadi *et al.*, 2019].

### 2.1 Hierarchical Reinforcement Learning

The standard RL formulation considers an agent selecting primitive actions at every time-step. In hierarchical RL [Barto and Mahadevan, 2003], the agent instead selects temporally extended actions or *skills*. These skills are commonly modeled as options [Sutton *et al.*, 1999], where each option  $o$  is described using three elements: the initiation region,  $\mathcal{I}_o : \mathcal{S} \rightarrow \{0, 1\}$ , describing the set of states from which the option can be executed; the termination region,  $\beta_o : \mathcal{S} \rightarrow \{0, 1\}$ , describing the subgoal region in which option execution must terminate; and the policy,  $\pi_o : \mathcal{S} \rightarrow \mathcal{A}$ , which drives the agent from  $\mathcal{I}_o$  to  $\beta_o$ . Several methods have been proposed to autonomously discover useful options (see Abel [2020], chapter 2.3, for a survey).

**Skill Chaining:** For two options  $o_i$  and  $o_j$  to be reliably sequentially executable, it must be the case that  $\beta_{o_i} \subseteq \mathcal{I}_{o_j}$ . The skill-chaining algorithm [Konidaris and Barto, 2009] explicitly satisfies this property by learning options such that  $\beta_{o_i} = \mathcal{I}_{o_j}$ . In skill-chaining, each discovered option  $o$  learns its initiation region  $\mathcal{I}_{o,\theta}$  using a binary classifier  $\theta$  that describes the region from which the  $\pi_o$  reaches  $\beta_o$  with high probability. Simultaneously,  $\pi_o$  is learned using a model-free RL algorithm to reach its subgoal region  $\beta_o$ . The algorithm is recursive: it first learns an option that initiates near the goal and reliably takes the agent to the goal; then it learns another option whose termination region is the initiation region of the first option; then it repeats the procedure targeting the new option. Options are chained together in this fashion until the agent discovers an option whose initiation region contains the start state. Deep skill chaining (DSC) [Bagaria and Konidaris, 2020] extended skill-chaining with deep RL, outperforming existing state-of-the-art skill-discovery algorithms [Levy *et al.*, 2019; Bacon *et al.*, 2017].

Since DSC *explicitly* constructs composable options, we build on top of it to show that our proposed augmentations can substantially improve the reliability of the resulting options. However, since any skill-discovery algorithm must eventually compose learned skills to solve goal-directed tasks, our insights could, at least in principle, improve the robustness of their solutions.

### 2.2 Related Work

**Robustness in HRL:** Several variants of the Option-Critic architecture [Bacon *et al.*, 2017] have showcased robustness to changes in transition dynamics or the reward function [Mankowitz *et al.*, 2018; Khetarpal and Precup, 2019; Tiwari and Thomas, 2019; Jain *et al.*, 2021; Harutyunyan *et al.*, 2019]. By contrast, we seek to improve the reliability of hierarchical methods in stationary, sparse-reward MDPs.

**Composing Options:** Some recent work has studied the problem of composing skills as linear [Barreto *et al.*, 2019]

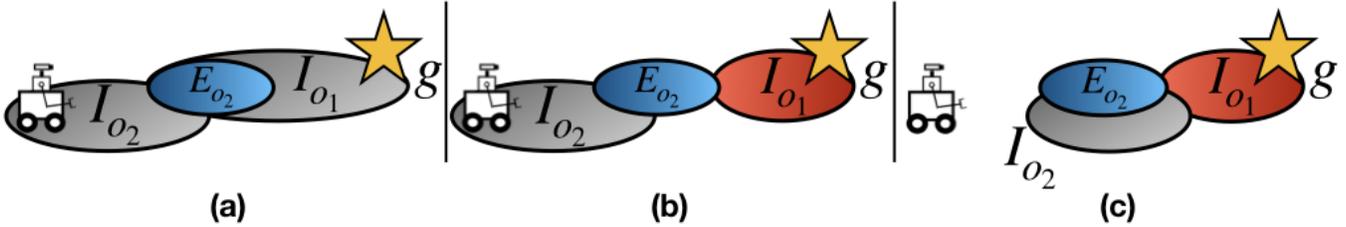


Figure 2: The non-stationary subgoal problem: (a) an agent has two options  $o_1, o_2$  such that  $o_2$  targets  $I_{o_1}$  and  $o_1$  targets goal  $g$ ;  $E_{o_2}$  is the next state distribution of  $\pi_{o_2}$ . (b) If  $I_{o_1}$  shrinks, executing  $o_2$  no longer allows the agent to execute  $o_1$  (since  $E_{o_2} \not\subseteq I_{o_1}$ ). (c) This causes  $\beta_{o_2}$  to shift forward, which invalidates the previous policy  $\pi_{o_2}$ , which in turn causes  $I_{o_2}$  to shrink.

or non-linear [Qureshi *et al.*, 2020] combinations of options available to the agent. These methods require pre-trained options, while we discover them from scratch.

**Model-Based Skill Discovery:** Empowerment-driven skill-discovery methods [Gregor *et al.*, 2016; Eysenbach *et al.*, 2019] have recently been augmented with model-based RL via the DADS algorithm [Sharma *et al.*, 2020]. This led to a family of related methods that optimize for exploration [Campos Camuñez *et al.*, 2020], lifelong learning [Lu *et al.*, 2021], and skill acquisition in image-based observation spaces [Baumli *et al.*, 2020]. These methods represent substantial progress in unsupervised skill-discovery, but our setting differs from theirs in a few ways. First, unlike these methods, we do not require a special pre-training phase for skill-discovery. Second, while DADS is designed for the multi-task setting, we focus on creating more robust solutions to single-goal MDPs. Finally, skills learned by DADS have global support, whereas we consider composability for skills that specialize in different regions.

### 3 Robustly Learning Composable Options

The core difficulty in learning composable options is that, during learning, all three components ( $\mathcal{I}_o, \beta_o, \pi_o$ ) of every option are simultaneously in flux. This difficulty is compounded by the relationship between the initiation region of one option and the subgoal region of another—changes to one option’s initiation region changes another’s subgoal, in turn changing its own policy and initiation region. These changes cascade to downstream options, propagating instability and causing chains of composable options to become unreliable or even break. This situation—which can happen whenever composable options are being learned simultaneously—is illustrated in Figure 2.

To ameliorate this difficulty, which we call *nonstationary subgoals*, we propose to robustify the learning process of all three option components. First, we propose a two-classifier representation of the initiation region that provides a stable subgoal target while encouraging exploration. Second, we use hindsight-experience replay to learn a generalized policy that enables us to target subgoal states that maximize the probability of being able to execute the successor option. Finally, to manage the difficulty of learning so many functions in parallel, we leverage the higher stability of model-based RL methods to learn more robust option policies.

#### 3.1 Dual Initiation Classifiers to Avoid Shrinkage

Previous approaches learned a single binary classifier to represent  $\mathcal{I}_o$  for each option  $o$  in the skill chain. If the option execution succeeds (i.e, the agent reaches  $\beta_o$ ), it adds a new positive example for further refining  $\mathcal{I}_o$ . Since the positive example is from a region already inside the classifier, a successful execution leaves the classifier’s decision boundary largely unchanged. Alternatively, if the option execution fails to reach  $\beta_o$ , the agent gets a negative example for the next training iteration of  $\mathcal{I}_o$ . Since this negative example comes from a region inside the classifier, it often shrinks the initiation region over time, with no opportunity to expand.

To avoid this issue, we propose a dual-classifier parameterization of  $\mathcal{I}_o$ —representing it using both optimistic and pessimistic classifiers. The pessimistic classifier represents the states from which we are highly confident that option execution will succeed, and so is a stable region for other options to target. However, if the agent could only ever choose to execute the option from inside this classifier, exploration would be hindered because the option would be prevented from expanding to new regions. To encourage exploration outside the pessimistic region, we also use an optimistic classifier to represent states where the agent can choose to execute the option. Eventually, the two classifiers should converge to approximate the “true” initiation region of the option.

Many techniques could be used to learn these two classifiers; we learn the optimistic classifier using a two-class SVM [Cortes and Vapnik, 1995] and the pessimistic classifier using a one-class SVM [Tax and Duin, 1999]. For more details, please refer to Section 3 of the Appendix.

#### 3.2 Robust Subgoals via Goal State Selection

We next turn to termination regions. The “chainability” of options  $o_i$  and  $o_j$  implies that the subgoal region  $\beta_{o_i}$  is a subset of the initiation region  $\mathcal{I}_{o_j}$ . As  $\mathcal{I}_{o_j}$  is refined over time, the subgoal region  $\beta_{o_i}$  also changes, and consequently, so does the option’s terminating reward. Since learning is highly sensitive to even small changes in the reward function [Packer *et al.*, 2018], this creates instability in learning  $\pi_{o_i}$  [Lu *et al.*, 2019]. DSC [Bagaria and Konidaris, 2020] mitigates this issue by freezing  $\mathcal{I}_{o_j}$  after a fixed number of learning iterations, which can lead to the agent being stuck with poor estimates of  $\mathcal{I}_{o_j}$  that hinder reliable skill composition.

Nevertheless, an option’s subgoal region will unavoidably change continually as its target option refines its initiation

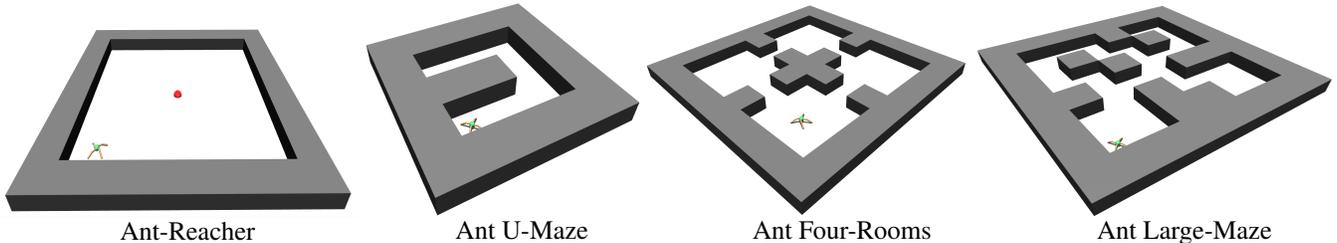


Figure 3: Maze navigation problems used to test our algorithm. These tasks require that the agent simultaneously learn good gait policies that stabilize the “ant” robot and navigate to a distant goal in the presence of unknown obstacles.

classifier. To be robust to changes in  $\beta_o$ , we propose to make each option policy more flexible: rather than a fixed policy  $\pi_o$ , each option learns a goal-conditioned policy  $\pi_o(s|g)$  using hindsight experience replay (HER). By conditioning the option policy  $\pi_o$  on goals sampled from  $\beta_o$ , and postponing selecting  $g$  until option execution time, the agent learns a policy robust to non-stationarity in  $\beta_o$ .

The goal-conditioned option policy strategy necessitates a strategy for sampling subgoal states from  $\beta_o$ . We propose optimizing two objectives for choosing subgoal states: (a) robustness and (b) hierarchical optimality.

**Selecting Subgoal States for Robustness.** Each option policy is rewarded for reaching its termination region; from its own perspective, all states in its termination region are equally rewarding. However, for a subsequent option  $o$ , some start states are better than others. How can we pick a subgoal for one option so that we increase the probability that a successive option execution will succeed?

One way is to evaluate the probability of every positive example of  $o$  succeeding (by querying the probabilistic classifier representing  $I_o$ ), but that becomes computationally very expensive as the agent gathers experience. Instead, we use the simple heuristic that, over time, the agent will learn to execute the option from reliable start states, and simply sample from the set of positive examples used to train  $\mathcal{I}_o$ .

**Selecting Subgoal States for Hierarchical Optimality.**

The method above results in feasible trajectories that are, at best, *recursively optimal* [Barto and Mahadevan, 2003]. We would prefer to pick subgoals for each option in the skill-chain so that the overall solution trajectory is approximately *hierarchically optimal* [Barto and Mahadevan, 2003].

To select such subgoals, we first store the states  $\hat{\beta}_o$  in which each option  $o$  triggered  $\beta_o$ . Then, we use dynamic programming (DP) to distribute the value of reaching the MDP’s goal  $g$  to all the  $\hat{\beta}_o$ s along the skill chain. This results in a value table  $\hat{Q} : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  that can be used to pick a subgoal  $s_g$  for the current option  $o_t$  from state  $s_t$ :  $s_g = \arg \max_{s \in \hat{\beta}_{o_t}} \hat{Q}(s_t, s)$ .

The quality of the resulting sub-goals depends on how well  $\hat{\beta}_o$  approximates  $\beta_o$ . If it is a perfect approximation, this DP algorithm yields the hierarchically optimal solution; otherwise, it yields a near-optimal solution. For more details about the DP algorithm, including the pseudocode and a discussion about the algorithm’s time-complexity, please refer to Section 1 of the Appendix.

### 3.3 Learning Robust Option Policies

Finally, we consider the third component of an option: its policy. Given the number of components simultaneously being learned in hierarchical algorithms, policy learning must be highly stable for the agent to succeed. Although model-free methods are commonplace for learning option policies, model-based methods are often more stable and sample-efficient [Deisenroth and Rasmussen, 2011].

We therefore propose learning option policies using model-based RL. We follow Nagabandi *et al.* [2018] to learn a dynamics model  $f_\xi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ . However, Equation 1 is insufficient for good action-selection in sparse reward problems. As a result, we solve the following infinite-horizon optimization problem and then execute the first action [Lowrey *et al.*, 2019]:

$$\pi(s|g) = \arg \max_{a \in \mathcal{U}^{|\mathcal{A}|}} \sum_{t=1}^H \gamma^{t-1} \mathcal{R}(s_t, a) + \gamma^H V_\phi(s_H|g). \quad (2)$$

We follow the same procedure as the model-free variant of our algorithm to learn the terminal value function  $V_\phi$ .

## 4 Experiments

Our experiments aim to answer the following questions: 1) do the proposed improvements increase the probability with which a sequence of options can be successfully composed? 2) Does the subgoal selection algorithm approximate hierarchically optimal trajectories? 3) How does the proposed algorithm compare to flat RL and other skill-discovery algorithms?

To answer these questions, we use a test-bed comprising four continuous-control maze-navigation tasks (shown in Figure 3) involving an “ant” robot simulated using MuJoCo [Todorov *et al.*, 2012; Duan *et al.*, 2016; Fu *et al.*, 2020].

**Sparse vs Dense Rewards.** Dense reward functions, although commonplace in RL, are often problematic because they demand cumbersome engineering [Yu *et al.*, 2020], can lead to sub-optimal solutions [Ng *et al.*, 1999] or reduce the problem so much that simple search might outperform RL [Mania *et al.*, 2018]. Since hierarchies are a promising way to address the challenges of sparse rewards, we evaluate all algorithms in the sparse reward setting.

### 4.1 Implementation Details

We use TD3 [Fujimoto *et al.*, 2018] and HER to learn goal-conditioned value functions in all our experiments. Transi-

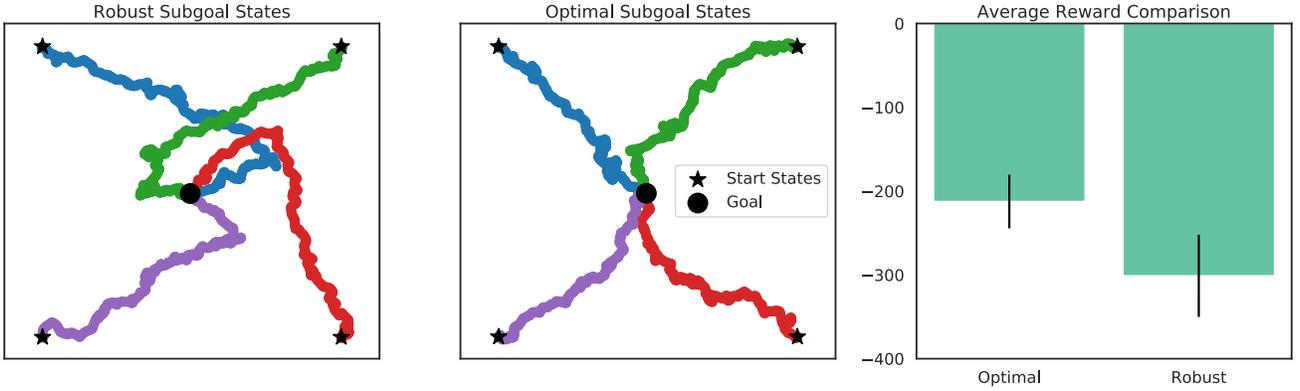


Figure 4: DSC trajectories (only the coordinates of the CoM are visualized) in Ant-Reacher when selecting subgoals (*left*) for robustness and (*middle*) using our dynamic programming algorithm. (*right*) When using DP, the agent scores better average reward (averaged over 5 runs; bars represent standard error; higher is better).

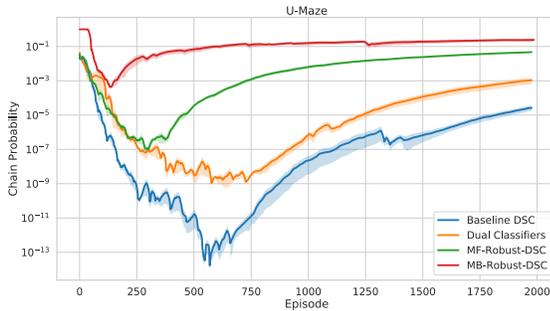


Figure 5: Comparing the robustness of skill chains discovered in the Ant U-Maze domain. For a description of the “chain probability” metric, please refer to Section 4.2. Vertical axis in log-scale; curves are averaged over 5 runs; shaded regions denote standard error.

tions seen by all options were used to train a single dynamics model  $f_\xi$ , critic  $V_\phi$ , and actor  $\pi_\zeta$ .  $V_\phi$  and  $\pi_\zeta$  were used for selecting actions in the model-free case. Only  $V_\phi$  from TD3 was used in the model-based case, where we approximately solved Equation 2 for action-selection. For more details, see Section 4 of the Appendix.

## 4.2 Evaluating Robustness

We first evaluate whether the proposed changes increase the composability of discovered options, by measuring the robustness of skill-chains constructed in Ant U-Maze.

We measure the robustness of a sequence of options as the probability that executing each option would take the agent to a state from which it could successfully execute the next option (until it finally reaches the goal). We call this the *chain-ing probability*:

$$p_{chain}(o_1, \dots, o_N) = \prod_{i=1}^N p(s_i \sim \beta_{o_{i-1}} \in \mathcal{I}_{o_i}), \quad (3)$$

where  $p(s_i \sim \beta_{o_{i-1}} \in \mathcal{I}_{o_i})$  represents the open-loop probability that the agent will be able to execute  $o_i$  after executing

option  $o_{i-1}$ . We approximate each probability term by the empirical “success rate” of the option:

$$p(s_i \sim \beta_{o_{i-1}} \in \mathcal{I}_{o_i}) \approx \frac{\#successes(o_{i-1})}{\#executions(o_{i-1})},$$

where  $\#successes(o_{i-1})$  is the number of times option  $\pi_{o_{i-1}}$  was able to reach  $\beta_{o_{i-1}}$ .

We ablate our proposals by comparing the robustness of the following versions of our algorithm:

1. Baseline DSC: The deep skill chaining algorithm from Bagaria and Konidaris [2020] that we aim to improve.
2. Dual Classifiers: We add the dual classifier approach from Section 3.1 to the baseline DSC.
3. MF-Robust-DSC: We add goal-conditioned policies to (2); this is the model-free version of the algorithm from Section 3.2.
4. MB-Robust-DSC: We add model-based policies to (3); this is the version of the algorithm from Section 3.3.

**Results.** Figure 5 shows that each of our proposals successively increases the robustness of the discovered skill-chain. The shape of the curves warrants discussion: as the agent discovers new options, its chain probability, at first, decreases. This is for two reasons: (a) more terms between 0 and 1 are included in the product of Equation 3 and (b) when initialized, option policies are not strong enough to reach their subgoals. Over time, two factors are responsible for the increasing trend: (a) the agent has discovered as many options as it needs to solve the problem, at which point no more terms are added to the product and (b) each option’s policy improves. Finally, the small absolute values on the vertical axis is due to the metric being an open-loop probability—we report the closed-loop success rate of the algorithm later in Section 4.4.

## 4.3 Evaluating Hierarchical Optimality

To compare the two subgoal selection strategies outlined in section 3.2, we ran the model-based variant of our algorithm on the Ant-Reacher domain. The goal state corresponds to the center of mass (CoM) of the ant being at (0, 0); the starting

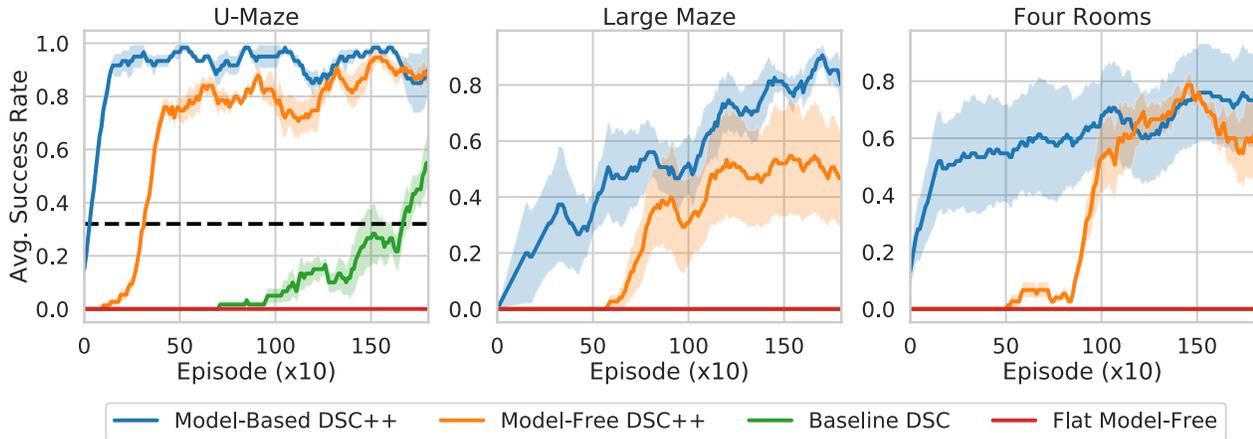


Figure 6: Average success rate in sparse-reward ant-navigation problems. Shaded regions denote standard error (averaged over 5 runs). The black dashed line in the leftmost subplot is the success rate of the flat model-free baseline after 12,000 episodes of training. In the middle and rightmost subplots, the green line is hidden below the red one.

position of the ant is sampled uniformly from  $(-10, 10)$ . The learned initiation classifiers tended to lie in approximately concentric circles around the goal state (visualized in Figure 1 of the Appendix). After both algorithm variants were trained for 1000 episodes, they were evaluated on how many steps it would take them to reach the goal when starting in the four corners of the domain  $\{(\pm 9, \pm 9)\}$ .

**Results.** Robust subgoal selection often yielded subgoal states on the “other side” of the goal—if the robot was on the bottom-left of the domain, it could pick a subgoal closer to the top-right. This led to the sub-optimal trajectories shown in Figure 4a, where the ant would often over-shoot the goal. By contrast, the dynamic programming procedure picked subgoal states that were always between the agent’s current state and the overall goal—leading to much more sensible, and approximately hierarchically optimal, solution trajectories (Figure 4b). Figure 4c shows that the agent using our DP algorithm for picking subgoals collected higher average rewards, further suggesting smoother overall trajectories to the goal.

#### 4.4 Learning Curves

Our final experiment compares the following algorithms on ant U-maze, large-maze and four-rooms:

1. Flat model-free baseline: we used TD3+HER [Fujimoto *et al.*, 2018; Andrychowicz *et al.*, 2017] as a flat model-free baseline, since it is a state-of-the-art method for continuous control (we used the same algorithm to learn our model-free option policies).
2. Flat finite-horizon model-based baseline: we used the model-based algorithm from Nagabandi *et al.* [2018] as our flat model-based baseline (we used the same algorithm to learn our model-based option policies).
3. Flat infinite-horizon model-based baseline: given that our problems are sparse-reward, we augment the model-based baseline with TD3, which is used to learn a value function that informs action-selection (Equation 2).

4. Deep skill chaining (DSC): we used DSC as our HRL baseline, because we extend DSC and it outperformed other skill-discovery methods [Bagaria and Konidaris, 2020].
5. Model-Free DSC++ (ours): the model-free variant of our algorithm described in Section 3.2.
6. Model-based DSC++ (ours): the model-based variant of our algorithm described in Section 3.3.

We report the “average success rate” metric from Andrychowicz *et al.* [2017]. Every 10 episodes, we ran the algorithm from a fixed start state  $(0, 0)$ , checking if it could reach the goal within 1000 steps. Due to its simplicity, we use the robust subgoal selection algorithm from Section 3.3 for rolling out our goal-conditioned option policies.

**Results.** Figure 6 shows the average success rate of competing methods averaged over 5 random seeds. Both versions of the flat model-based baseline were unable to achieve a  $> 0\%$  success rate; so we leave them out of the learning curves in Figure 6. Ant U-Maze was the only sparse-reward problem that TD3 and the baseline DSC were able to achieve a  $> 0\%$  (but TD3 had to be trained for far more episodes). The proposed algorithm (shown with blue and orange curves) were easily able to outperform all baselines including DSC.

## 5 Conclusion

The success of skill discovery in goal-directed tasks depends on learning options that can be sequentially composed. However, robust composition is challenging because of non-stationary subgoal regions. We proposed methods that address all three components of option learning—initiation regions, termination conditions, and option policies—to learn reliably composable options. We experimentally showed that our augmentations measurably improve the robustness of discovered skills and eventually allow us to solve more challenging, long-horizon problems.

## Acknowledgements

We thank Kshitij Sachan, Cam Allen, Sam Lobel, and other members of the Brown BigAI group for their suggestions. This research was supported by NSF grants 1955361, 1717569 and the DARPA Lifelong Learning Machines program under grant FA8750-18-2-0117. This research was conducted using computational resources and services at the Center for Computation and Visualization, Brown University.

## References

- D Abel. *A Theory of Abstraction in Reinforcement Learning*. PhD thesis, Brown University, 2020.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hind-sight experience replay. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Kavosh Asadi, Dipendra Misra, Seungchan Kim, and Michel L Littman. Combating the compounding-error problem with a multi-step model. *arXiv preprint arXiv:1905.13320*, 2019.
- Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Akhil Bagaria and George Konidaris. Option discovery using deep skill chaining. In *International Conference on Learning Representations*, 2020.
- André Barreto, Diana Borsa, Shaobo Hou, Gheorghe Comanici, Eser Aygün, Philippe Hamel, Daniel Toyama, Shibl Mourad, David Silver, and Doina Precup. The option keyboard: Combining skills in reinforcement learning. In *Advances in Neural Information Processing Systems*, 2019.
- Andrew G Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 2003.
- Kate Baumli, David Warde-Farley, Steven Hansen, and Volodymyr Mnih. Relative variational intrinsic control. *arXiv preprint arXiv:2012.07827*, 2020.
- Matthew M Botvinick, Yael Niv, and Andrew G Barto. Hierarchically organized behavior and its neural foundations: a reinforcement learning perspective. *Cognition*, 113(3), 2009.
- Robert R Burridge, Alfred A Rizzi, and Daniel E Koditschek. Sequential composition of dynamically dexterous robot behaviors. *The International Journal of Robotics Research*, 1999.
- Víctor Campos Camúñez, Alex Trott, Caiming Xiong, Richard Socher, Xavier Giró Nieto, and Jordi Torres Viñals. Explore, discover and learn: unsupervised discovery of state-covering skills. In *International Conference on Machine Learning*, 2020.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 1995.
- Marc Deisenroth and Carl E Rasmussen. PILCO: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning*, 2011.
- Thomas G Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. *JAIR*, 2000.
- Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*, 2016.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.
- Kevin Frans, Jonathan Ho, Xi Chen, Pieter Abbeel, and John Schulman. Meta learning shared hierarchies. In *International Conference on Learning Representations*, 2018.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, 2018.
- Carlos E Garcia, David M Prett, and Manfred Morari. Model predictive control: theory and practice—a survey. *Automatica*, 25(3):335–348, 1989.
- Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *ArXiv*, abs/1611.07507, 2016.
- Anna Harutyunyan, Will Dabney, Diana Borsa, Nicolas Heess, Remi Munos, and Doina Precup. The termination critic. In *The 22nd International Conference on Artificial Intelligence and Statistics*, 2019.
- Arushi Jain, Khimya Khetarpal, and Doina Precup. Safe option-critic: learning safety in the option-critic architecture. *The Knowledge Engineering Review*, 36, 2021.
- Yuu Jinnai, Jee Won Park, David Abel, and George Konidaris. Discovering options for exploration by minimizing cover time. In *International Conference on Machine Learning*, 2019.
- Leslie Pack Kaelbling and Tomás Lozano-Pérez. Learning composable models of parameterized skills. In *IEEE International Conference on Robotics and Automation*, 2017.
- Leslie Pack Kaelbling. Learning to achieve goals. In *International Joint Conference on Artificial Intelligence*, 1993.
- Khimya Khetarpal and Doina Precup. Learning options with interest functions. In *AAAI*, 2019.
- George Konidaris and Andrew Barto. Building portable options: Skill transfer in reinforcement learning. In *International Joint Conference on Artificial Intelligence*, 2007.
- George Konidaris and Andrew Barto. Skill discovery in continuous reinforcement learning domains using skill chaining. *Advances in Neural Information Processing Systems*, 22, 2009.

- George Konidaris, Leslie Pack Kaelbling, and Tomas Lozano-Perez. From skills to symbols: Learning symbolic representations for abstract high-level planning. *JAIR*, 2018.
- Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Hierarchical reinforcement learning with hindsight. In *International Conference on Learning Representations*, 2019.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Long-Ji Lin. Reinforcement learning for robots using neural networks. Technical report, Carnegie-Mellon Univ Pittsburgh PA School of Computer Science, 1993.
- Kendall Lowrey, Aravind Rajeswaran, Sham Kakade, Emanuel Todorov, and Igor Mordatch. Plan online, learn offline: Efficient learning and exploration via model-based control. In *International Conference on Learning Representations*, 2019.
- Tomas Lozano-Perez, Matthew T Mason, and Russell H Taylor. Automatic synthesis of fine-motion strategies for robots. *The International Journal of Robotics Research*, 1984.
- Kevin Lu, Igor Mordatch, and Pieter Abbeel. Adaptive online planning for continual lifelong learning. *arXiv preprint arXiv:1912.01188*, 2019.
- Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Reset-free lifelong learning with skill-space planning. In *International Conference on Learning Representations*, 2021.
- Aleksandr Mikhailovich Lyapunov. The general problem of the stability of motion. *International journal of control*, 1992.
- Horia Mania, Aurelia Guy, and Benjamin Recht. Simple random search of static linear policies is competitive for reinforcement learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 1805–1814, 2018.
- Daniel Mankowitz, Timothy Mann, Pierre-Luc Bacon, Doina Precup, and Shie Mannor. Learning robust options. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 2015.
- Anusha Nagabandi, Gregory Kahn, Ronald S Fearing, and Sergey Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *IEEE International Conference on Robotics and Automation*, 2018.
- Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, volume 99, 1999.
- Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song. Assessing generalization in deep reinforcement learning. *arXiv preprint arXiv:1810.12282*, 2018.
- Ahmed H Qureshi, Jacob J Johnson, Yuzhe Qin, Taylor Henderson, Byron Boots, and Michael C Yip. Composing task-agnostic policies with deep reinforcement learning. *International Conference on Learning Representations*, 2020.
- Jette Randløv, Andrew G Barto, and Michael T Rosenstein. Combining reinforcement learning with a local control algorithm. In *International Conference on Machine Learning*, 2000.
- Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In *International Conference on Machine Learning*, 2015.
- Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020.
- Farzaneh Shoeleh and Masoud Asadpour. Graph based skill acquisition and transfer learning for continuous reinforcement learning domains. *Pattern Recognition Letters*, 2017.
- Avi Singh, Albert Yu, Jonathan Yang, Jesse Zhang, Aviral Kumar, and Sergey Levine. COG: Connecting new skills to past experience with offline reinforcement learning. In *Conference on Robot Learning*, 2020.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- R.S. Sutton, , D. Precup, and S. Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 1999.
- David MJ Tax and Robert PW Duin. Support vector domain description. *Pattern recognition letters*, 1999.
- Russ Tedrake. LQR-trees: Feedback motion planning on sparse randomized trees. 2009.
- Saket Tiwari and Philip S Thomas. Natural option critic. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012.
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 1992.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Metaworld: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, 2020.

# Appendix: Robustly Learning Composable Options in Deep Reinforcement Learning

Akhil Bagaria\*, Jason Senthil\*, Matthew Slivinski and George Konidaris

Department of Computer Science, Brown University

{akhil\_bagaria, jason\_senthil, matthew\_slivinski}@brown.edu, gdk@cs.brown.edu

## 1 Picking Hierarchically Optimal Subgoals

In this section, we provide more details about the algorithm discussed in Section 3.2 of the main paper. Specifically, we pick (approximately) hierarchically optimal subgoal states for each option in the chain using a dynamic programming algorithm. By approximately “hierarchically optimal”, we mean that the resulting trajectory approximates the optimal solution to the overall problem, *given acquired options* [Barto and Mahadevan, 2003].

---

### Algorithm 1 Subgoal selection algorithm

---

```

1: procedure CREATE-VALUE-TABLE( $o_1, o_2, \dots, o_N$ )
2:   for  $o_i \in (o_1, \dots, o_N)$  do  $\triangleright$  Start from the goal-option
3:     update-option-values( $o_i$ )
4:   end for
5: end procedure
6: procedure UPDATE-OPTION-VALUES( $o$ )
7:   for  $s \in \mathcal{I}_o$  do  $\triangleright$  Samples from initiation set
8:     for  $s' \in \mathcal{E}_o$  do  $\triangleright$  Samples from termination set
9:       if  $o$  is not the goal-option then
10:         $p = o.\text{parent}$ 
11:         $\tilde{Q}(s, s') = V_o(s|s') + \gamma \max_{s'' \in \mathcal{E}_p} \tilde{Q}(s', s'')$ 
12:       else
13:         $\tilde{Q}(s, s') = \mathcal{R}(s, s')$ 
14:       end if
15:     end for
16:   end for
17: end procedure
18: procedure PICK-SUBGOAL( $s_t, o_t$ )
19:   if  $o_t$  is the goal-option then
20:     return  $\mathcal{M}.\text{goal-state}$ 
21:   end if
22:    $p = o.\text{parent}$ 
23:   return  $\arg \max_{s \in \mathcal{I}_p} \tilde{Q}(s_t, s)$ 
24: end procedure

```

---

The algorithm operates on pairs of input, output states for each option  $o$ . The input states are states inside  $\mathcal{I}_o$  and the output states  $\mathcal{E}_o$  are the subset of samples from the option’s termination set that still satisfy its parent option’s pessimistic

classifier (i.e, they are samples from the option’s current termination region).

**Algorithm Description.** The goal state of the first option in the chain (the “goal-option”) is known—it is the goal state  $g$  of the MDP  $\mathcal{M}$ . As a result, the values of all the states in  $\mathcal{I}_o$ , when  $o$  is the goal-option is the terminal reward corresponding to reaching  $g$  (line 13 of Algorithm 1). Next, we examine the input-output states of the next option—the one that targets the “goal-option”. Each pair of input-output states is assigned a value based on that option’s value-function—*bootstrapped* by the parent option’s value (line 11 of Algorithm 1). This process is recursively continued until the input-output states of each option in the chain have been assigned a value in the table  $\tilde{Q}$ . The `create-value-table()` procedure is called at the end of every episode.

**Time complexity.** Given that the agent has  $K$  options, and the option  $o$  with the most states in its effect set has  $|\mathcal{E}_o| = N$ , then the big-o complexity of this algorithm is  $\mathcal{O}(K \times N^2)$ .

Once the value table  $\tilde{Q}$  is constructed, the agent first picks the option to execute:  $o_t = \pi_{\mathcal{O}}(s_t)$ . Then, it queries the value table  $\tilde{Q}$  to find the next subgoal state from the current state-option pair  $(s_t, o_t)$ :

$$s_g = \arg \max_{s \in \mathcal{E}_{o_t}} \tilde{Q}(s_t, s).$$

## 2 Full Algorithm Pseudo-Code

Algorithms 2 and 3 outline the pseudo-code of our method.

## 3 Initiation Set classifiers

In this section, we provide more details about how the initiation classifiers are trained and utilized. Then, we will present some qualitative comparison with the classifiers learned by the baseline DSC algorithm.

**Training Initiation Classifiers.** We represent each option’s initiation region using an optimistic classifier and a pessimistic classifier. We parameterized the optimistic classifier using a two-class SVM  $\theta_1$  which was trained using positive  $X^+$  (states from which option execution was successful) and negative  $X^-$  (states from which option execution failed) examples. We then pass  $X^+$  through  $\theta_1$  and the subset of it which is still predicted as positive,  $X^{++}$ , is used to train the one-class SVM  $\theta_2$ . This leads to a tight pessimistic classifier

---

\*Equal contribution

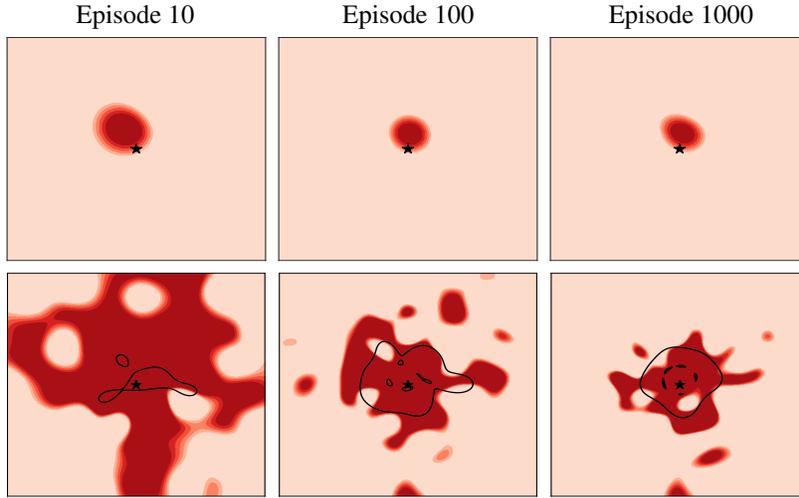


Figure 1: Initiation regions learned for the option that drives the agent to the goal state (visualized as a star in the center of each plot) in Ant-Reacher. The horizontal and vertical axes represent the  $x, y$  positions of the ant’s CoM. The top row visualizes the classifiers learned by the baseline DSC algorithm; the bottom row corresponds to the Robust DSC algorithm. In the bottom row, the optimistic classifier is visualized in dark red, the pessimistic classifier as the black contour lines. The Robust DSC agent learns a bigger initiation region, which implies larger goal regions for successive options, thereby making their policy learning process easier.

---

#### Algorithm 2 Robust DSC Rollout

---

```

1: procedure ROLLOUT(Skill Chain  $\mathcal{O}$ , Option Horizon  $H$ )
2:   Initialize empty trajectory buffer  $\mathcal{B}$ 
3:   for each timestep  $t$  do
4:     Select option  $o$  using policy over options  $\pi_{\mathcal{O}}(s_t)$ 
5:     Sample a goal for selected option:  $g \sim \beta_o$ 
6:     Execute option policy  $\pi_o(\cdot | g)$  in the environment
7:     Add trajectory  $\tau = \bigcup_{i=0}^{H-1} (s_i, o, a_i, s_{i+1}, g)$  to  $\mathcal{B}$ 
8:     if final state  $s_H$  reached goal  $g$  then
9:       Add  $\tau$  to  $o$ ’s list of positive examples
10:    else
11:      Add  $\tau$  to  $o$ ’s list of negative examples
12:    end if
13:    Refit option  $o$ ’s initiation classifier
14:    Add  $\tau$  to  $o$ ’s replay buffer and update  $\pi_o$ 
15:  end for
16:  return  $\mathcal{B} = \bigcup_t (s_t, o_t, a_t, s_{t+1}, g_t)$ 
17: end procedure

```

---

around the states from which we are highly confident about executing that option.

Both types of SVMs were trained using an RBF kernel [Scholkopf *et al.*, 1997]. We used the ThunderSVM software package [Wen *et al.*, 2018] to implement the initiation classifiers. Following related work [Levy *et al.*, 2019; Bagaria and Konidaris, 2020; Eysenbach *et al.*, 2019; Sharma *et al.*, 2020], the classifiers were trained on the subset of the state (ant position) that matters to the reward function.

**Querying Initiation Classifiers.** If option  $o_2$  targets option  $o_1$ , a state  $s$  is in  $\beta_{o_2}$  if  $s$  is inside  $o_1$ ’s pessimistic classifier. To make the optimistic classifier larger than the pessimistic one, option  $o_1$  can be executed if  $s$  is inside *either* classifier’s positive region.

---

#### Algorithm 3 Robust DSC Algorithm

---

```

1: procedure ROBUST-DSC(Start state  $s_0$ , Goal region  $g$ )
2:   Initialize global option  $o_G$  such that  $\mathcal{I}_{o_G}(\cdot) = 1$ 
3:   Initialize goal option  $o_g$  such that  $\beta_{o_g} = g$ 
4:   Initialize skill chain  $\mathcal{O}$  with  $\{o_g\}$ 
5:   for each episode do
6:     transitions = ROLLOUT( $\mathcal{O}, H$ )
7:     if  $s_0 \notin \mathcal{I}_o, \forall o \in \mathcal{O}$  then
8:       Create new option  $\omega$ 
9:       Add  $\omega$  to skill chain  $\mathcal{O}$ 
10:    end if
11:  end for
12: end procedure

```

---

**Qualitative Evaluation of Learned Classifiers.** We now qualitatively evaluate whether our proposed changes can discover better initiation regions than baseline DSC. We compare the two algorithms in the Ant-Reacher domain where the ant has to navigate to the center of  $20 \times 20$  arena. We will visualize the initiation classifier for the “goal option”, which is the option that targets the MDP’s goal state (shown with a star in figure 1).

Figure 1 shows the initiation classifiers learned by the baseline DSC algorithm (top) and the Robust DSC algorithm (bottom). As is visually clear, the baseline algorithm discovers an overly conservative initiation classifier which shrinks over-time; implying a smaller goal region for a successive option that would target it. By contrast, using the new algorithm, the optimistic classifier (shown as the red region in the bottom row) shrinks and the pessimistic classifier (shown with black contours) expands until the two converge on a intuitively reasonable estimate of the option’s initiation region.

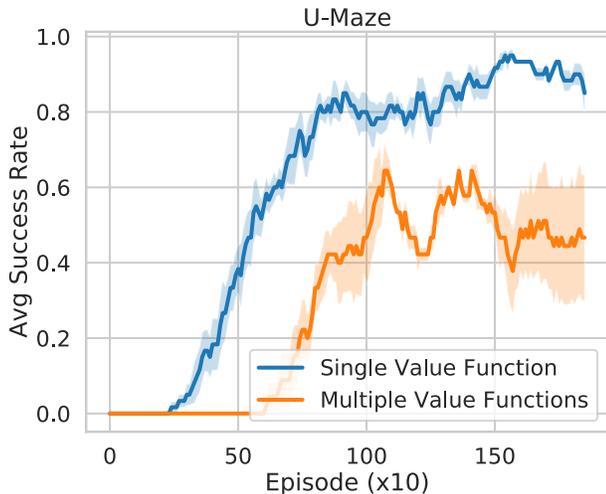


Figure 2: Learning curves comparing two versions of the model-free variant of our algorithm in Ant U-Maze. In one, every option distills its experiences into the same value function (blue); in the other, each option learns a separate, independent value function (orange).

## 4 Policy learning implementation details

Each option  $o$ 's reward function  $\mathcal{R}_o$  gives a reward of 0 for reaching its subgoal and  $-1$  otherwise:  $\mathcal{R}_o(s, g \in \beta_o) = \mathcal{R}(s, g)$ . We then use an off-the-shelf actor-critic algorithm, TD3 [Fujimoto *et al.*, 2018] and HER [Andrychowicz *et al.*, 2017] for policy learning. In the model-based setting, we found that using TD3 with output normalization [Agarwal *et al.*, 2020] was more stable than a vanilla TD3, where the data collected for learning  $V_\phi$  is “more off-policy” [Levine *et al.*, 2020].

**Learning a single value function.** Usually each option has to learn a different value function that is tied to its own subgoal region. However, due to goal-conditioning, we can now re-use the same value function to pick actions for different options. Each option conditions the value function on goal states drawn from its own subgoal region, and then picks an action accordingly:  $\pi_o(s|g \sim \beta_o) = \arg \max_{a \in \mathcal{A}} Q_\phi(s, a|g)$ . We hypothesize that learning a single value function would be better because it distills data collected by all options into a single function [Rusu *et al.*, 2015], precluding the need to re-learn similar representations over and over. We ablate this design decision of using a single value function rather than having a separate one for each option—the results are presented in Figure 2. This parameterization of options suggests connections between HRL and curriculum learning [Narvekar *et al.*, 2020].

### 4.1 Deep Model-based RL

We adopt the approach from Nagabandi *et al.* [2018] for our learned dynamics model. We parameterize our learned dynamics function  $f_\xi(s_t, a_t)$  as a feedforward deep neural network which take as input the current state  $s_t$  and action  $a_t$ , and predicts the change in state  $s_t$  over the time step duration of  $\Delta t$ . Thus the next state  $\hat{s}_{t+1}$  is predicted by  $\hat{s}_{t+1} = s_t + f_\xi(s_t, a_t)$ .

## Learning the Dynamics Model

The neural network architecture for our dynamics model has two hidden layers, each with dimension 500, with LeakyReLU activations. We train this model using the Adam optimizer with learning rate 0.001 and batch size 1024. Again following Nagabandi *et al.* [2018], our model-based experiments begin with 50 episodes of random rollouts from  $(0, 0)$ . This data is used to train the first iteration of the dynamics model  $f_\xi$ , which is later fine-tuned during RL. After the first iteration of training, we fit  $f_\xi$  at the end of every episode. Unlike Nagabandi *et al.* [2018], we do not add noise to the state during training or inference.

## Model-Predictive Control

We utilize our dynamics model by using model predictive control. We follow Nagabandi *et al.* [2018] identically here. We use a random-shooting method in which  $M$  action sequences of length  $K$  are randomly generated from the uniform random distribution  $\mathcal{U}(-1, 1)$ —the first action in the cost minimizing sequence is executed in the environment. In all of our experiments,  $M = 14,000$  and  $K = 7$ .

## 5 Hyperparameter Settings

For TD3, we use the hyperparameters exactly as in Fujimoto *et al.* [2018], except for the learning rate value in both actor and critic networks. We found that a lower learning rate was helpful for stabilizing value function learning. The hyperparameters related to TD3 are listed in Tables 1 and 2; those related to training  $f_\xi$  are in Table 3. Finally, the option specific hyperparameters are listed in Table 4.

Parameter	Value
Replay buffer size	1e6
Critic Learning rate	$10^{-5}$
Actor Learning rate	$10^{-5}$
Optimizer	Adam
Target Update Rate $\tau$	$5 \cdot 10^{-3}$
Batch size	100
Iterations per time step	1
Discount Factor	0.99
Output Normalization	True

Table 1: TD3 Hyperparameters for Model-based DSC++

## 6 Number of Skills over Time

The skill-chaining algorithm creates as many options as it needs to reliably solve a goal-directed task. As a result, it creates more options to solve harder tasks than easier ones. Figure 3 shows how the number of skills change over time for different configurations of the DSC algorithm. Baseline DSC requires significantly more options to solve the same problem with the same reliability as the DSC++ algorithm. In addition to the results presented in Figure 4 of the main paper, this suggests that our proposals improve the reliability of the skill-chain.

Parameter	Value
Replay buffer size	1e6
Critic Learning rate	$3 \cdot 10^{-4}$
Actor Learning rate	$3 \cdot 10^{-4}$
Optimizer	Adam
Target Update Rate $\tau$	$5 \cdot 10^{-3}$
Batch size	100
Iterations per time step	1
Discount Factor	0.99
Output Normalization	False

Table 2: TD3 Hyperparameters for Model-free DSC++

Parameter	Value
Batch size	1024
Optimizer	Adam
Controller horizon $H$	7
Number actions sampled $K$	14000

Table 3: Hyperparameters for learning dynamics model  $f_\xi$

## 7 Compute Infrastructure

All experiments in this paper were run on the Brown University CCV compute cluster which included 8 NVIDIA QuadroRTX GPUs and 16 Intel Skylake CPUs. Running our experiments requires only 1 GPU at a time, the multiple GPUs were used to run the same code with different random seeds in parallel.

## References

- Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*, 2020.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hind-sight experience replay. In *Advances in Neural Information Processing Systems*, 2017.
- Akhil Bagaria and George Konidaris. Option discovery using deep skill chaining. In *ICLR*, 2020.
- Andrew G Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 13(1-2):41–77, 2003.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and

Parameter	U-Maze	Large Maze	Four Rooms
Gestation Period	5	5	10
Option Timeout	200	200	200
Buffer Length	50	50	50

Table 4: Skill Chaining Hyperparameters

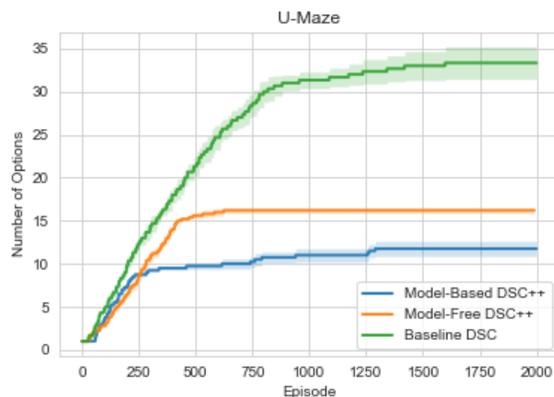


Figure 3: Number of discovered skills over time on Ant U-Maze.

Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *ICLR*, 2019.

Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1582–1591, 2018.

Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.

Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Hierarchical reinforcement learning with hindsight. In *International Conference on Learning Representations*, 2019.

Anusha Nagabandi, Gregory Kahn, Ronald S Fearing, and Sergey Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7559–7566. IEEE, 2018.

Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E Taylor, and Peter Stone. Curriculum learning for reinforcement learning domains: A framework and survey. *Journal of Machine Learning Research*, 21(181):1–50, 2020.

Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy distillation. *arXiv preprint arXiv:1511.06295*, 2015.

Bernhard Scholkopf, Kah-Kay Sung, Christopher JC Burges, Federico Girosi, Partha Niyogi, Tomaso Poggio, and Vladimir Vapnik. Comparing support vector machines with gaussian kernels to radial basis function classifiers. *IEEE transactions on Signal Processing*, 45(11):2758–2765, 1997.

Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised dis-

covery of skills. In *International Conference on Learning Representations (ICLR)*, 2020.

Zeyi Wen, Jiashuai Shi, Qinbin Li, Bingsheng He, and Jian Chen. ThunderSVM: A fast SVM library on GPUs and CPUs. *Journal of Machine Learning Research*, 19:797–801, 2018.