
Robust Option Learning for Adversarial Generalization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Compositional reinforcement learning is a promising approach for training poli-
2 cies to perform complex long-horizon tasks. Typically, a high-level task is decom-
3 posed into a sequence of subtasks and a separate policy is trained to perform each
4 subtask. In this paper, we focus on the problem of training subtask policies in a
5 way that they can be used to perform any task; here, a task is given by a sequence
6 of subtasks. We aim to maximize the worst-case performance over all tasks as
7 opposed to the average-case performance. We formulate the problem as a two
8 agent zero-sum game in which the adversary picks the sequence of subtasks. We
9 propose two RL algorithms to solve this game: one is an adaptation of existing
10 multi-agent RL algorithms to our setting and the other is an asynchronous version
11 which enables parallel training of subtask policies. We evaluate our approach on
12 two multi-task environments with continuous states and actions and demonstrate
13 that our algorithms outperform state-of-the-art baselines.

14 1 Introduction

15 Reinforcement learning (RL) has proven to be a promising strategy for solving complex control
16 tasks such as walking [13], autonomous driving [17], and dexterous manipulation [3]. However, a
17 key challenge facing the deployment of reinforcement learning in real-world tasks is its high sample
18 complexity—to solve any new task requires training a new policy designed to solve that task. One
19 promising solution is *compositional reinforcement learning*, where individual *options* (or *skills*) are
20 first trained to solve simple tasks; then, these options can be composed together to solve more
21 complex tasks [25, 24, 17]. For example, if a driving robot learns how to make left and right turns
22 and to drive in a straight line, it can then drive along any route composed of these primitives.

23 A key challenge facing compositional reinforcement learning is the generalizability of the learned
24 options. In particular, options trained under one distribution of tasks may no longer work well if used
25 in a new task, since the distribution of initial states from which the options are used may shift. An
26 alternate approach is to train the options separately to perform specific subtasks, but options trained
27 this way might cause the system to reach states from which future subtasks are hard to perform. One
28 can overcome this issue by handcrafting rewards to encourage avoiding such states [17], in which
29 case they generalize well, but this approach relies heavily on human time and expertise.

30 We propose a novel framework that addresses this challenge by formulating the option learning
31 problem as an adversarial reinforcement learning problem. At a high level, the adversary chooses
32 the task that minimizes the reward achieved by composing the available options. Thus, the goal is
33 to learn a set of *robust options* that perform well across *all* potential tasks. Then, we provide two
34 algorithms for solving this problem. The first adapts existing ideas for using reinforcement learning
35 to solve Markov games to our setting. Then, the second shows how to leverage the compositional

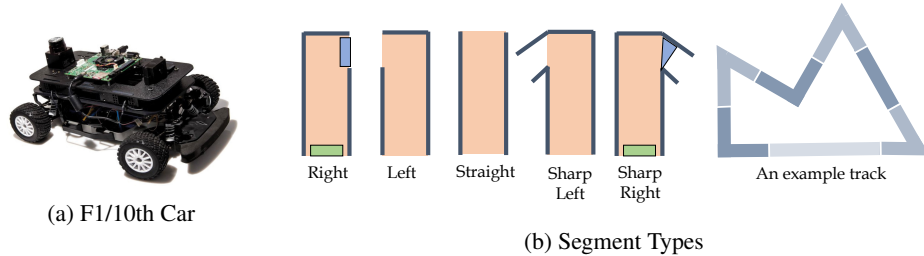


Figure 1: F1/10th Environment. The entry and exit regions for the right and sharp right segments are shown in green and blue respectively.

36 structure of our problem to learn options in parallel at each step of a value iteration procedure; in
 37 some cases, by enabling such parallelism, we can reduce the computational cost of learning.

38 We validate our approach on two benchmarks: (i) a rooms environment where a point mass robot
 39 must navigate any given sequence of rooms, where the sequence is an arbitrary combination of
 40 straight, left, and right turns, and (ii) a simulated version of the F1/10th car, where a small racing car
 41 must navigate any racetrack composed of several different track segments. In both, our empirical
 42 results demonstrate that robust options are critical for performing well on a wide variety of tasks.

43 In summary, our contributions are: (i) a game theoretic formulation of the compositional reinforcement
 44 learning problem, (ii) two algorithms for solving this problem, and (iii) an empirical evaluation
 45 demonstrating the effectiveness of our approach.

46 **Motivating example.** Let us consider a small scale autonomous racing car shown in Figure 1 (a).
 47 We would like to train a controller that can be used to navigate the car through *all* tracks constructed
 48 using five kinds of segments; the possible segments are shown in Figure 1 (b) along with an example
 49 track. The state of the car is a vector (x, y, v, θ) where (x, y) is its position on the track relative to
 50 the current segment, v is its current speed and θ is the heading angle. An action is a pair $(a, \omega) \in \mathbb{R}^2$
 51 where a is the throttle input and ω is the steering angle. In this environment, completing each
 52 segment is considered a subtask and a task corresponds to completing a sequence of segments—
 53 e.g., *straight* \rightarrow *right* \rightarrow *left* \rightarrow *sharp-right*. Upon completion of a subtask, the car enters
 54 the next segment and a change-of-coordinates is applied to the car’s state which is now relative to the
 55 new segment. The goal here is to learn one option for each subtask such that the agent can perform
 56 any task using these options.

57 If one trains the options independently with the only goal of reaching the end of each segment
 58 (e.g., using distance-based rewards), it might (and does) happen that the car reaches the end of a
 59 segment in a state that was not part of the initial states used to train the policy corresponding to the
 60 next subtask. Therefore, one should make sure that the initial state distribution used during training
 61 includes such states as well—either manually or using dataset aggregation [38]. Furthermore, it is
 62 possible that the car reaches a state in the exit region of a segment from which it is challenging to
 63 complete the next subtask—e.g., a state in which the car is close to and facing towards a wall. Our
 64 algorithm identifies during training that, in order to perform future subtasks, it is better to reach the
 65 end of a segment in a configuration where the car is facing straight relative to the next segment. As
 66 demonstrated in our experiments, this leads to robust options and improved sample efficiency.

67 **Related work.** The options framework [41] is commonly used to model subtask policies as tem-
 68 porally extended actions. In hierarchical RL [32, 31, 22, 9, 5, 43], options are trained along with a
 69 high-level controller that chooses the sequence of options to execute in order to complete a specific
 70 high-level task. There is also work on discovering options—e.g., using intrinsic motivation [30],
 71 entropy maximization [10], semi-supervised RL [12], skill chaining [20], expectation maximiza-
 72 tion [8] and subgoal identification [40]. There has also been a lot of research on planning using
 73 learned options [1, 18, 37, 42, 21].

74 There has been some work on RL for zero-shot generalization [44, 33, 39, 23, 4]; however, in prior
 75 work, the learning objective is to maximize average performance with respect to a fixed distribution
 76 over tasks as opposed to the worst-case. Some hierarchical RL algorithms have also been shown to
 77 enable few-shot generalization [18] to unseen tasks. Most closely related to our work is the work

78 on compositional RL in the multi-task setting [17] in which the subtask policies are trained using
 79 standard RL algorithms in a naive way without guarantees regarding worst-case performance.

80 There has also been work on skill composition using transition policies [25]; this method assumes
 81 that the subtask policies are fixed and learns one transition policy per subtask which takes the system
 82 from an end state to a “favourable” initial state for the subtask. However, poorly trained subtask
 83 policies can lead to situations in which it is not possible to achieve such transitions. In contrast, our
 84 approach trains subtask policies which compose well without requiring additional transition policies.
 85 A recent paper [24] proposes a framework for training subtask policies with the aim of composing
 86 them to perform a complex long-horizon task. However, their approach assumes that the high-level
 87 task is fixed and the options are trained to maximize the performance with respect to a specific task.

88 There has been a lot of research on multi-agent RL algorithms [29, 15, 16, 28, 35, 36, 2] including
 89 algorithms for two-agent zero-sum games [6, 45, 27]. In this paper, we utilize the specific structure
 90 of our game to obtain a simple algorithm that neither requires solving matrix games nor trains a
 91 separate policy for the adversary. Furthermore, we show that we can obtain an asynchronous RL
 92 algorithm which enables learning options in parallel.

93 2 Problem Formulation

94 A *multi-task Markov decision process (MDP)* is a tuple $\mathcal{M} = (S, A, P, \Sigma, R, F, T, \gamma, \eta, \sigma_0)$, where
 95 S are the states, A are the actions, $P(s' | s, a) \in [0, 1]$ is the probability of transitioning from s to
 96 s' on action a , η is the initial state distribution, and $\gamma \in (0, 1)$ is the discount factor. Furthermore,
 97 Σ is a set of subtasks and for each subtask $\sigma \in \Sigma$, $R_\sigma : S \times A \rightarrow \mathbb{R}$ is a reward function¹, $F_\sigma \subseteq S$
 98 is a set of final states where the subtask is considered completed and $T_\sigma : F_\sigma \times S \rightarrow [0, 1]$ is the
 99 jump probability function; upon reaching a state s in F_σ the system jumps to a new state s' with
 100 probability $T_\sigma(s' | s)$. For the sake of clarity, we assume² that $T_\sigma(s' | s) = 0$ for any s' with
 101 $s' \in F_{\sigma'}$ for some σ' . Finally, $\sigma_0 \in \Sigma$ is the initial subtask which has to be completed first³. A
 102 multi-task MDP can be viewed as a discrete time variant of a hybrid automaton model [17].

103 In the case of our motivating example, the set of subtasks is given by

$$\Sigma = \{\text{left}, \text{right}, \text{straight}, \text{sharp-left}, \text{sharp-right}\}$$

104 with F_σ denoting the exit region of the segment corresponding to subtask σ . We use the jump
 105 transitions T to model the change-of-coordinates performed upon reaching an exit region. The
 106 reward function R_σ for a subtask σ is given by $R_\sigma(s, a, s') = -\|s' - c_\sigma\|_2^2 + B \cdot \mathbb{1}(s' \in F_\sigma)$ where
 107 c_σ is the center of the exit region and the subtask completion bonus B is a positive constant.

A task τ is defined to be an infinite sequence⁴ of subtasks $\tau = \sigma_0\sigma_1\dots$, and \mathcal{T} denotes the set of all
 tasks. For any task $\tau \in \mathcal{T}$, $\tau[i]$ denotes the i^{th} subtask σ_i in τ . In our setting, the task is chosen by
 the environment nondeterministically. Given a task τ , a configuration of the environment is a pair
 $(s, i) \in S \times \mathbb{Z}_{\geq 0}$ with $s \notin F_{\tau[i]}$ denoting that the system is in state s and the current subtask is $\tau[i]$.
 The initial distribution over configurations $\Delta : S \times \mathbb{Z}_{\geq 0} \rightarrow [0, 1]$ is given by $\Delta(s, i) = \eta_{\tau[0]}(s)$ if
 $i = 0$ and 0 otherwise. The probability of transitioning from (s, i) to (s', j) on an action a is

$$\Pr((s', j) | (s, i), a) = \begin{cases} \sum_{s'' \in F_{\tau[i]}} P(s'' | s, a) T_{\tau[i]}(s' | s'') & \text{if } j = i + 1 \\ P(s' | s, a) & \text{if } j = i \\ 0 & \text{otherwise.} \end{cases}$$

Intuitively, the system transitions to the next subtask if the current subtask is completed and stays in
 the current subtask otherwise. A (deterministic) policy is a function $\pi : S \rightarrow A$, where $a = \pi(s)$
 is the action to take in state s . Our goal is to learn one policy π_σ for each subtask σ such that the
 discounted reward over the worst-case task τ is maximized. Formally, given a set of policies $\Pi =$
 $\{\pi_\sigma | \sigma \in \Sigma\}$ and a task τ , we can define a Markov chain over configurations with initial distribution
 Δ and transition probabilities given by $P_\Pi((s', j) | (s, i)) = \Pr((s, j') | (s, i), \pi_{\tau[i]}(s))$. We denote

¹We can also have $R_\sigma : S \times A \times S \rightarrow \mathbb{R}$ depending on the next state but we omit it for clarity of presentation.

²This assumption can be removed by adding a fictitious copy of F_σ to S for each $\sigma \in \Sigma$.

³When there is no fixed initial subtask, we can add a fictitious initial subtask.

⁴A finite sequence can be appended with an infinite sequence of a fictitious subtask with zero reward.

by \mathcal{D}_τ^Π the distribution over infinite sequences of configurations $\rho = (s_0, i_0)(s_1, i_1) \dots$ generated by τ and Π . Then, we define the objective function as

$$J(\Pi) = \inf_{\tau \in \mathcal{T}} \mathbb{E}_{\rho \sim \mathcal{D}_\tau^\Pi} \left[\sum_{t=0}^{\infty} \gamma^t R_{\tau[i_t]}(s_t, \pi_{\tau[i_t]}(s_t)) \right].$$

108 These definitions can be naturally extended to stochastic policies as well. In our motivating ex-
 109 ample, choosing a large enough completion bonus B guarantees the discounted reward to be
 110 higher for runs in which more subtasks are completed. Our aim is to compute a set of policies
 111 $\Pi^* \in \arg \max_{\Pi} J(\Pi)$. Each subtask policy π_σ defines an option [41] $o_\sigma = (\pi_\sigma, I_\sigma, \beta_\sigma)$ where
 112 $I_\sigma = S \setminus F_\sigma$ and $\beta_\sigma(s) = \mathbb{1}(s \in F_\sigma)$. Here, the choice of which option to trigger is made by the
 113 environment rather than the agent.

114 3 Reduction to Stagewise Markov Games

The problem statement naturally leads to a game theoretic view in which the environment is the adversary. We can formally reduce the problem to a two-agent zero-sum Markov game $\mathcal{G} = (\bar{S}, A_1, A_2, \bar{P}, \bar{R}, \bar{\gamma}, \bar{\eta})$ where $\bar{S} = S \times \Sigma$ is the set of states, $A_1 = A$ are the actions of agent 1 (the agent learning the options) and $A_2 = \Sigma$ are the actions of agent 2 (the adversary). The transition probability function $\bar{P} : \bar{S} \times A_1 \times A_2 \times \bar{S} \rightarrow [0, 1]$ is given by

$$\bar{P}((s', \sigma') | (s, \sigma), a_1, a_2) = \begin{cases} P(s' | s, a_1) & \text{if } s \notin F_\sigma \text{ \& } \sigma = \sigma' \\ T_\sigma(s' | s) & \text{if } s \in F_\sigma \text{ \& } \sigma' = a_2 \\ 0 & \text{otherwise.} \end{cases}$$

115 We observe that the states are partitioned into two sets $\bar{S} = S_1 \cup S_2$ where $S_1 = \{(s, \sigma) | s \notin F_\sigma\}$
 116 is the set of states where agent 1 acts (causing a step in \mathcal{M}) and $S_2 = \{(s, \sigma) | s \in F_\sigma\}$ is the set
 117 of states where agent 2 takes actions (causing a change of subtask); this makes \mathcal{G} a stagewise game.
 118 The reward function $\bar{R} : \bar{S} \times A_1 \rightarrow \mathbb{R}$ is given by $\bar{R}((s, \sigma), a) = R_\sigma(s, a)$ if $s \notin F_\sigma$ and 0 otherwise.
 119 The discount factor depends on the state and is given by $\bar{\gamma}(s, \sigma) = \gamma$ if $s \notin F_\sigma$ and 1 otherwise; this
 120 is because a change of subtask does not invoke a step in \mathcal{M} . The initial state distribution $\bar{\eta}$ is given
 121 by $\bar{\eta}(s, \sigma) = \eta(s) \mathbb{1}(\sigma = \sigma_0)$. A run of the game is a sequence $\bar{\rho} = \bar{s}_0 a_0^1 a_0^2 \bar{s}_1 a_1^1 a_1^2 \dots$ where $\bar{s}_t \in \bar{S}$
 122 and $a_t^i \in A_i$.

A (deterministic) policy for agent i is a function $\pi_i : \bar{S} \rightarrow A_i$. Given policies π_1 and π_2 for agents 1 and 2, respectively and a state $\bar{s} \in \bar{S}$ we denote by $\mathcal{D}_{\bar{s}}^{\mathcal{G}}(\pi_1, \pi_2)$ the distribution over runs generated by π_1 and π_2 starting at \bar{s} . Then, the value of a state \bar{s} is defined by

$$V^{\pi_1, \pi_2}(\bar{s}) = \mathbb{E}_{\bar{\rho} \sim \mathcal{D}_{\bar{s}}^{\mathcal{G}}(\pi_1, \pi_2)} \left[\sum_{t=0}^{\infty} \left(\prod_{k=0}^{t-1} \bar{\gamma}(\bar{s}_k) \right) \bar{R}(\bar{s}_t, a_t^1) \right].$$

We are interested in computing a policy π_1^* maximizing

$$J_{\mathcal{G}}(\pi_1) = \mathbb{E}_{\bar{s} \sim \bar{\eta}} \left[\min_{\pi_2} V^{\pi_1, \pi_2}(\bar{s}) \right].$$

123 Given a policy π_1 for agent 1, we can construct a policy π_σ for any subtask σ given by $\pi_\sigma(s) =$
 124 $\pi_1(s, \sigma)$; we denote by $\Pi(\pi_1)$ the set of subtask policies constructed this way. The following the-
 125 orem connects the objective of the game with our multi-task learning objective; all proofs are in
 126 Appendix A.

127 **Theorem 3.1.** *For any policy π_1 for agent 1 in \mathcal{G} , we have $J(\Pi(\pi_1)) \geq J_{\mathcal{G}}(\pi_1)$.*

128 Therefore, $J_{\mathcal{G}}(\pi_1)$ is a lower bound on the objective $J(\Pi(\pi_1))$ which we seek to maximize. Now,
 129 let us define the optimal value of a state \bar{s} by $V^*(\bar{s}) = \max_{\pi_1} \min_{\pi_2} V^{\pi_1, \pi_2}(\bar{s})$. The following
 130 theorem shows that it is possible to construct a policy π_1^* that maximizes $J_{\mathcal{G}}(\pi_1)$ from the optimal
 131 value function V^* .

Theorem 3.2. *For any policy π_1^* such that for all $(s, \sigma) \in S_1$,*

$$\pi_1^*(s, \sigma) \in \arg \max_{a \in A} \left\{ \bar{R}((s, \sigma), a) + \gamma \cdot \sum_{s' \in S} P(s' | s, a) V^*(s', \sigma) \right\},$$

132 *we have that $\pi_1^* \in \arg \max_{\pi_1} J_{\mathcal{G}}(\pi_1)$.*

Algorithm 1 Asynchronous value iteration algorithm for computing optimal subtask policies.

```

1: function ASYNCDVALUEITERATION( $\mathcal{M}, V$ )
2:   while stopping criterion is met do
3:     for  $\sigma \in \Sigma$  do // in parallel
4:       Compute  $\mathcal{W}_\sigma(V)$ 
5:        $V \leftarrow \mathcal{F}_{\text{async}}(V)$  // using Equation 3

```

133 3.1 Value Iteration

134 In this section, we briefly look at two value iteration algorithms to compute V^* which we later adapt
135 in Section 4 to obtain learning algorithms. Let $\mathcal{V} = \{V : S_1 \rightarrow \mathbb{R}\}$ denote the set of all value
136 functions over S_1 . Given a value function $V \in \mathcal{V}$ we define its extension to all of \bar{S} using

$$\llbracket V \rrbracket(s, \sigma) = \begin{cases} \min_{\sigma' \in \Sigma} \sum_{s' \in S} T_\sigma(s' | s) V(s', \sigma') & \text{if } s \in F_\sigma \\ V(s, \sigma) & \text{otherwise.} \end{cases} \quad (1)$$

137 For a state $s \in F_\sigma$, $\llbracket V \rrbracket(s, \sigma)$ denotes the worst-case value (according to V) with respect to the
138 possible choices of next subtask σ' . Now, we consider the Bellman operator $\mathcal{F} : \mathcal{V} \rightarrow \mathcal{V}$ defined by

$$\mathcal{F}(V)(s, \sigma) = \max_{a \in A} \left\{ \bar{R}((s, \sigma), a) + \gamma \cdot \sum_{s' \in S} P(s' | s, a) \llbracket V \rrbracket(s', \sigma) \right\} \quad (2)$$

139 for all $(s, \sigma) \in S_1$. Let us denote by $V^* \downarrow_{S_1}$ the restriction of V^* to S_1 . The following lemma
140 follows straightforwardly giving us our first value iteration procedure.

141 **Theorem 3.3.** \mathcal{F} is a contraction mapping with respect to the ℓ_∞ -norm on \mathcal{V} and $V^* \downarrow_{S_1}$ is the
142 unique fixed point of \mathcal{F} with $\lim_{n \rightarrow \infty} \mathcal{F}^n(V) = V^* \downarrow_{S_1}$ for all $V \in \mathcal{V}$.

143 Next we consider an *asynchronous* value iteration procedure which allows us to parallelize com-
144 puting subtask policies for different subtasks. Given a subtask σ and a value function $V \in \mathcal{V}$, we
145 define a *subtask MDP* \mathcal{M}_σ^V which behaves like \mathcal{M} until reaching a final state $s \in F_\sigma$ after which it
146 transitions to a dead state \perp achieving a reward of $\llbracket V \rrbracket(s, \sigma)$. Formally, $\mathcal{M}_\sigma^V = (S_\sigma, A, P_\sigma, R_\sigma^V, \gamma)$
147 where $S_\sigma = S \sqcup \{\perp\}$ with \perp being a special dead state, $P_\sigma(s' | s, a) = P(s' | s, a)$ if $\perp \neq s \notin F_\sigma$
148 and $P_\sigma(s' | s, a) = \mathbb{1}(s' = \perp)$ otherwise. The reward function is given by $R_\sigma^V(s, a) = R_\sigma(s, a)$ if
149 $\perp \neq s \notin F_\sigma$, $R_\sigma^V(s, a) = \llbracket V \rrbracket(s, \sigma)$ if $\perp \neq s \in F_\sigma$ and is 0 otherwise. We denote by $\mathcal{W}_\sigma(V)$ the
150 optimal value function of the MDP \mathcal{M}_σ^V . We then define the asynchronous value iteration operator
151 $\mathcal{F}_{\text{async}} : \mathcal{V} \rightarrow \mathcal{V}$ using

$$\mathcal{F}_{\text{async}}(V)(s, \sigma) = \mathcal{W}_\sigma(V)(s). \quad (3)$$

152 We can show that repeated application of $\mathcal{F}_{\text{async}}$ leads to the optimal value function V^* of the \mathcal{G} .

153 **Theorem 3.4.** For any $V \in \mathcal{V}$, $\lim_{n \rightarrow \infty} \mathcal{F}_{\text{async}}^n(V) \rightarrow V^* \downarrow_{S_1}$.

Since each $\mathcal{W}_\sigma(V)$ can be computed independently, we can parallelize the computation of $\mathcal{F}_{\text{async}}$
giving us the algorithm in Algorithm 1. We can also show that it is not necessary to compute $\mathcal{W}_\sigma(V)$
exactly. Let $\mathcal{V}_\sigma = \{\bar{V} : S_\sigma \rightarrow \mathbb{R}\}$ be the set of all value functions over S_σ . For a fixed $V \in \mathcal{V}$, let
 $\mathcal{F}_{\sigma, V} : \mathcal{V}_\sigma \rightarrow \mathcal{V}_\sigma$ denote the usual Bellman operator for the MDP \mathcal{M}_σ^V given by

$$\mathcal{F}_{\sigma, V}(\bar{V})(s) = \max_{a \in A} \left\{ R_\sigma^V(s, a) + \gamma \cdot \sum_{s' \in S_\sigma} P_\sigma(s' | s, a) \bar{V}(s') \right\}$$

for all $\bar{V} \in \mathcal{V}_\sigma$ and $s \in S_\sigma$. For any $V \in \mathcal{V}$ and $\sigma \in \Sigma$, we define a corresponding $V_\sigma \in \mathcal{V}_\sigma$ using
 $V_\sigma(s) = \llbracket V \rrbracket(s, \sigma)$ if $s \in S$ and $V_\sigma(\perp) = 0$. Then, for any integer $m > 0$ and $V \in \mathcal{V}$, we can use
 $\mathcal{F}_{\sigma, V}^m(V_\sigma)$ as an approximation to $\mathcal{W}_\sigma(V)$. Let us define $\mathcal{F}_m : \mathcal{V} \rightarrow \mathcal{V}$ using

$$\mathcal{F}_m(V)(s, \sigma) = \mathcal{F}_{\sigma, V}^m(V_\sigma)(s).$$

154 Intuitively, \mathcal{F}_m corresponds to performing m steps of value iteration in each subtask MDP \mathcal{M}_σ^V
155 (which can be parallelized) starting at V_σ . The following theorem guarantees convergence when
156 using \mathcal{F}_m instead of $\mathcal{F}_{\text{async}}$.

157 **Theorem 3.5.** For any $V \in \mathcal{V}$ and $m > 0$, $\lim_{n \rightarrow \infty} \mathcal{F}_m^n(V) \rightarrow V^* \downarrow_{S_1}$.

Algorithm 2 Robust Option Soft Actor Critic.

Inputs: Learning rates $\alpha_\psi, \alpha_\theta$, entropy weight β and Polyak rate δ .

```

1: function ROSAC( $\alpha_\psi, \alpha_\theta, \beta, \delta$ )
2:   Initialize parameters  $\{\psi_\sigma\}_{\sigma \in \Sigma}, \{\psi_\sigma^{\text{targ}}\}_{\sigma \in \Sigma}$  and  $\{\theta_\sigma\}_{\sigma \in \Sigma}$ 
3:   Initialize replay buffer  $\mathcal{B}$ 
4:   for each iteration do
5:     for each episode do
6:        $s_0 \sim \eta$ 
7:        $\sigma_0 \leftarrow \text{InitialSubtask}$ 
8:       for each step  $t$  do
9:          $a_t \sim \pi_{\theta_{\sigma_t}}(\cdot | s_t)$  and  $s_{t+1} \sim P(\cdot | s, a)$ 
10:         $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, s_{t+1})\}$ 
11:        if  $s_{t+1} \in F_{\sigma_t}$  then
12:           $s_{t+1} \sim T_{\sigma_t}(\cdot | s_{t+1})$ 
13:           $\sigma_{t+1} \leftarrow \text{Greedy}(\varepsilon, \arg \min_\sigma \tilde{V}(s_{t+1}, \sigma), \Sigma)$ 
14:        else
15:           $\sigma_{t+1} \leftarrow \sigma_t$ 
16:        for each gradient step do
17:          Sample batch  $B \sim \mathcal{B}$ 
18:          for  $\sigma \in \Sigma$  do
19:             $\psi_\sigma \leftarrow \psi_\sigma - \alpha_\psi \nabla_{\psi_\sigma} \mathcal{L}_Q(\psi_\sigma, B)$ 
20:             $\theta_\sigma \leftarrow \theta_\sigma - \alpha_\theta \nabla_{\theta_\sigma} \mathcal{L}_\pi(\theta_\sigma, B)$ 
21:             $\psi_\sigma^{\text{targ}} \leftarrow \delta \psi_\sigma + (1 - \delta) \psi_\sigma^{\text{targ}}$ 

```

158 **4 Learning Algorithms**

159 In this section, we present RL algorithms for solving the game \mathcal{G} . We first consider the finite MDP
160 setting for which we can obtain a modified Q -learning algorithm with a convergence guarantee. We
161 then present two algorithms based on Soft Actor Critic (SAC) for the continuous state setting.

162 **4.1 Finite MDP**

163 Assuming finite states and actions, we can obtain a Q -learning variant for solving \mathcal{G} which we call
164 *Robust Option Q -learning*. We assume that jump transitions T are known to the learner; this is usu-
165 ally the case since jump transitions are used to model subtask transitions and change-of-coordinates
166 within the controller. However, we believe that the algorithm can be easily extended to the scenario
167 where T is unknown.

168 We maintain a function $Q : S_1 \times A \rightarrow \mathbb{R}$ with $Q(s, \sigma, a)$ denoting $Q((s, \sigma), a)$. The corre-
169 sponding value function V_Q is defined using $V_Q(s, \sigma) = \max_{a \in A} Q(s, \sigma, a)$ and is extended to
170 all of \bar{S} as $\llbracket V_Q \rrbracket$. Note that, given a Q -function, the extended value function $\llbracket V_Q \rrbracket$ can be com-
171 puted exactly. Robust Option Q -learning is an iterative process—in each iteration t , it takes a step
172 $((s, \sigma), a_1, a_2, (s', \sigma))$ in \mathcal{G} with $(s, \sigma) \in S_1$ and updates the Q -function using

$$Q_{t+1}(s, \sigma, a_1) \leftarrow (1 - \alpha_t) Q_t(s, \sigma, a_1) + \alpha_t (\bar{R}((s, \sigma), a_1) + \gamma \llbracket V_{Q_t} \rrbracket(s', \sigma)). \quad (4)$$

173 where Q_t is the Q -function in iteration t and $\llbracket V_{Q_t} \rrbracket$ is the corresponding extended value function.

174 Under standard assumptions on the learning rates α_t , similar to Q -learning, we can show that Robust
175 Option Q -learning converges to the optimal Q -function almost surely. Here, the optimal Q -function
176 is defined by $Q^*(s, \sigma, a) = \bar{R}((s, \sigma), a) + \gamma \sum_{s' \in S} P(s' | s, a) V^*(s', \sigma)$ for all $(s, \sigma) \in S_1$. Let
177 $\alpha_t(s, \sigma, a)$ denote the learning rate used in iteration t if $Q_t(s, \sigma, a)$ is updated and 0 otherwise. Then,
178 we have the following theorem.

179 **Theorem 4.1.** *If $\sum_t \alpha_t(s, \sigma, a) = \infty$ and $\sum_t \alpha_t^2(s, \sigma, a) < \infty$ for all $(s, \sigma) \in S_1$ and $a \in A$, then*
180 $\lim_{t \rightarrow \infty} Q_t = Q^*$ *with probability 1.*

181 **4.2 Continuous States and Actions**

182 In the case of continuous states and actions, we can adapt any Q -function based RL algorithm such
 183 as Deep Deterministic Policy Gradients (DDPG) [26] or Soft Actor Critic (SAC) [14] to our setting.
 184 Here we present an SAC based algorithm that we call Robust Option SAC (ROSAC) which is outlined
 185 in Algorithm 2. This algorithm, like SAC, adds an entropy bonus to the reward function to improve
 186 exploration.

We maintain two Q -functions for each subtask σ , $Q_{\psi_\sigma} : S \rightarrow \mathbb{R}$ parameterized by ψ_σ and a target
 function $Q_{\psi_\sigma^{\text{targ}}}$ parameterized by $\psi_\sigma^{\text{targ}}$. We also maintain a stochastic subtask policy $\pi_{\theta_\sigma} : S \rightarrow$
 $\mathcal{D}(A)$ for each subtask σ where $\mathcal{D}(A)$ denotes the set of distributions over A . Given a step (s, a, s')
 in \mathcal{M} and a subtask σ with $s \notin F_\sigma$, we define the target value by

$$y_\sigma(s, a, s') = R_\sigma(s, a) + \gamma[V](s', \sigma)$$

where the value $[V](s', \sigma)$ is estimated using $\tilde{V}(s', \sigma) = Q_{\psi_\sigma^{\text{targ}}}(s', \tilde{a}) - \beta \log \pi_{\theta_\sigma}(\tilde{a} | s')$ with
 $\tilde{a} \sim \pi_{\theta_\sigma}(\cdot | s')$ if $s' \notin F_\sigma$. If $s' \in F_\sigma$, we estimate $[V](s', \sigma)$ using $\tilde{V}(s', \sigma) = \min_{\sigma' \in \Sigma} \tilde{V}(s'', \sigma')$
 where $\tilde{V}(s'', \sigma') = Q_{\psi_{\sigma'}^{\text{targ}}}(s'', \tilde{a}) - \beta \log \pi_{\theta_{\sigma'}}(\tilde{a} | s'')$ with $\tilde{a} \sim \pi_{\theta_{\sigma'}}(\cdot | s'')$ and $s'' \sim T_\sigma(\cdot | s')$.
 Now, given a batch $B = \{(s, a, s')\}$ of steps in \mathcal{M} we update ψ_σ using one step of gradient descent
 corresponding to the loss

$$\mathcal{L}_Q(\psi_\sigma, B) = \frac{1}{|B|} \sum_{(s, a, s') \in B} (Q_{\psi_\sigma}(s, a) - y_\sigma(s, a, s'))^2$$

and the subtask policy π_{θ_σ} is updated using the loss

$$\mathcal{L}_\pi(\theta_\sigma, B) = -\frac{1}{|B|} \sum_{(s, a, s') \in B} \mathbb{E}_{\tilde{a} \sim \pi_{\theta_\sigma}(\cdot | s)} [Q_{\psi_\sigma}(s, \tilde{a}) - \beta \log \pi_{\theta_\sigma}(\tilde{a} | s)].$$

187 The gradient $\nabla_{\theta_\sigma} \mathcal{L}_\pi(\theta_\sigma, B)$ can be estimated using the reparametrization trick if $\pi_{\theta_\sigma}(\cdot | s)$ is a
 188 Gaussian distribution whose parameters are differentiable w.r.t. θ_σ . We use Polyak averaging to
 189 update the target Q -networks $\{Q_{\psi_\sigma^{\text{targ}}} | \sigma \in \Sigma\}$.

190 Note that we do not train a separate policy for the adversary. During exploration, we use the ε -
 191 greedy strategy to select subtasks. We first estimate the ‘‘worst’’ subtask for a state s using $\tilde{\sigma} =$
 192 $\arg \min_{\sigma} \tilde{V}(s, \sigma)$ where $\tilde{V}(s, \sigma)$ is estimated as before. Then the function $\text{Greedy}(\varepsilon, \tilde{\sigma}, \Sigma)$ chooses
 193 $\tilde{\sigma}$ with probability $1 - \varepsilon$ and picks a subtask uniformly at random from Σ with probability ε .

194 **Asynchronous ROSAC.** We can also obtain an asynchronous version of the above algorithm which
 195 lets us train subtask policies in parallel. Asynchronous Robust Option SAC (AROSAC) is outlined in
 196 Algorithm 3. Here we use one replay buffer for each subtask. We maintain an initial state distribution
 197 $\tilde{\eta}$ over S to be used for training every subtask policy $\{\pi_\sigma\}_{\sigma \in \Sigma}$. $\tilde{\eta}$ is represented using a finite set of
 198 states D from which a state is sampled uniformly at random. The value function $\tilde{V} : S \times \Sigma \rightarrow \mathbb{R}$
 199 is estimated as before. To be specific, in each iteration, an estimate of any value $\tilde{V}(s, \sigma)$ is obtained
 200 on the fly using the Q -functions and the subtask policies from the previous iteration.

201 The SAC subroutine runs the standard Soft Actor Critic algorithm for N iterations on the subtask
 202 MDP $\mathcal{M}_\sigma^{\tilde{V}}$ (defined in Section 3)⁵ with initial state distribution $\tilde{\eta}$ (defaults to η if $D = \emptyset$). It returns
 203 the updated parameters along with states X_σ visited during exploration with $X_\sigma \subseteq F_\sigma$. The states
 204 in X_σ are used to update the initial state distribution for the next iteration following the Dataset
 205 Aggregation principle [38].

206 **5 Experiments**

207 We evaluate our algorithms ROSAC and AROSAC on two multi-task environments; a rooms environ-
 208 ment and an F1/10th racing car environment [11].

⁵Note that it is possible to obtain samples from $\mathcal{M}_\sigma^{\tilde{V}}$ as long as one can obtain samples from \mathcal{M} and
 membership in F_σ can be decided.

Algorithm 3 Asynchronous Robust Option Soft Actor Critic.

Inputs: Learning rates α , entropy weight β , Polyak rate δ and number of SAC iterations N .

```
1: function AROSAC( $\alpha, \beta, \delta, N$ )
2:   Initialize parameters  $\Psi = \{\psi_\sigma\}_{\sigma \in \Sigma}$ ,  $\Psi^{\text{targ}} = \{\psi_\sigma^{\text{targ}}\}_{\sigma \in \Sigma}$  and  $\Theta = \{\theta_\sigma\}_{\sigma \in \Sigma}$ 
3:   Initialize replay buffers  $\{\mathcal{B}_\sigma\}_{\sigma \in \Sigma}$  and Initialize  $D = \{\}$ 
4:   for each iteration do
5:      $\tilde{V} \leftarrow \text{OBTAINVALUEESTIMATOR}(\Psi, \Theta)$ 
6:     for  $\sigma \in \Sigma$  do // in parallel
7:        $\psi_\sigma, \psi_\sigma^{\text{targ}}, \theta_\sigma, X_\sigma \leftarrow \text{SAC}(\mathcal{M}_\sigma^{\tilde{V}}, D, \psi_\sigma, \psi_\sigma^{\text{targ}}, \theta_\sigma, \alpha, \beta, \delta, N)$ 
8:     for  $\sigma \in \Sigma$  do
9:       for  $s \in X_\sigma$  do
10:         $s' \sim T_\sigma(\cdot | s)$  and  $D \leftarrow D \cup \{s'\}$ 
```

209 **Rooms environment.** We consider the environment shown in Fig-
210 ure 2 which depicts a room with walls and exits. Initially the robot
211 is placed in the green triangle. The L-shaped obstacles indicate walls
212 within the room that the robot cannot cross. A state of the system is a
213 position $(x, y) \in \mathbb{R}^2$ and an action is a pair (v, θ) where v is the speed
214 and θ is the heading angle to follow during the next time-step. There
215 are three exits: left (blue), right (yellow) and up (grey) reaching each
216 of which is a subtask. Upon reaching an exit, the robot enters another
217 identical room where the exit is identified (via change-of-coordinates)
218 with the bottom entry region of the current room. A task is a sequence
219 of directions—e.g., left \rightarrow right \rightarrow up \rightarrow right indicating that
220 the robot should reach the left exit followed by the right exit in the
221 subsequent room and so on. Although the dynamics are simple, the
222 obstacles make learning challenging in the adversarial setting.

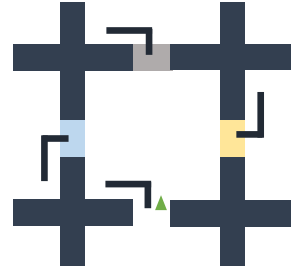


Figure 2: Rooms environment

223 **F1/10th environment.** This is the environment in the motivating example. A publicly available
224 simulator [11] of the F1/10th car is used for training and testing. The policies use the LiDAR
225 measurements from the car as input (as opposed to the state) and we assume that the controller can
226 detect the completion of each segment; as shown in prior work [17], one can train a separate neural
227 network to predict subtask completion.

228 **Baselines.** We compare our approach to three baselines. The baseline NAIVE trains one controller
229 for each subtask with the only aim of completing the subtask, similar to [17], using a manually
230 designed initial state distribution. DAGGER is a similar approach which, instead of using a manually
231 designed initial state distribution for training, infers the initial state distribution using the Dataset
232 Aggregation principle [38]. The MADDPG baseline solves the game \mathcal{G} using the multi-agent RL
233 algorithm proposed in [29] for solving concurrent Markov games with continuous states and actions.

234 **Evaluation.** We evaluate the performance of these algorithms against two adversaries. One adver-
235 sary is the random adversary which picks the next subtask uniformly at random from the set of all
236 subtasks. The other adversary estimates the worst sequence of subtasks for a given set of options
237 using Monte Carlo Tree Search (MCTS) [19]. The MCTS adversary is trained by assigning a reward
238 of 1 if it selects a subtask which the corresponding policy is unable to complete within a fixed time
239 budget and a reward of 0 otherwise. For the Rooms environment, we consider subtask sequences of
240 length at most 5 whereas for the F1/10th environment, we consider sequences of subtasks of length
241 at most 20. We evaluate both the average number of subtasks completed as well as the probability
242 of completing the set maximum number of subtasks.

243 **Results.** The plots for the rooms environment are shown in Figure 3 and plots for the F1/10th
244 environment are shown in Figure 4. We can observe that ROSAC is able outperform other approaches
245 and learn robust options. In the rooms environment, AROSAC achieves similar performance albeit
246 requiring more samples; however, it has the added benefit of being parallelizable. In the F1/10th
247 environment, it performs similar to the other baselines. DAGGER and NAIVE baselines are unable to
248 learn policies that can be used to perform long sequences of subtasks; this is mostly due to the fact

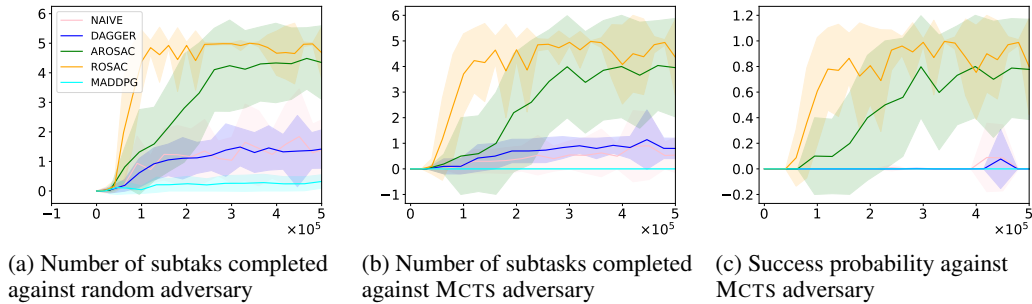


Figure 3: Plots for the Rooms environment. x -axis denoted the number of sample steps and y -axis denoted the either the average number of subtasks completed or the probability of completing 5 subtasks. Results are averaged over 10 runs. Error bars indicate \pm standard deviation.

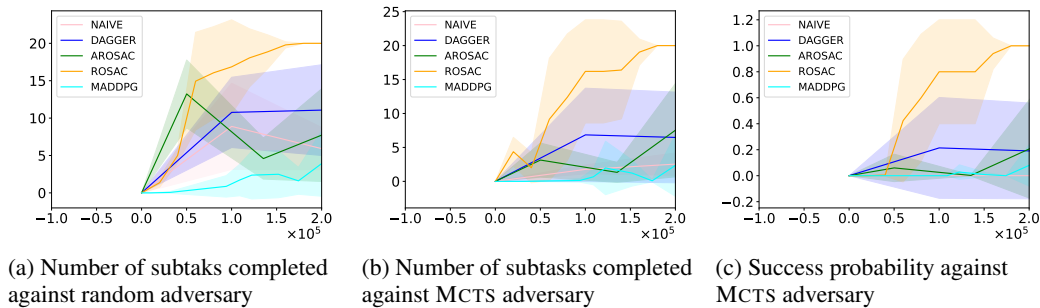


Figure 4: Plots for the F1/10th environment. x -axis denoted the number of sample steps and y -axis denoted the either the average number of subtasks completed or the probability of completing 20 subtasks. Results are averaged over 5 runs. Error bars indicate \pm standard deviation.

249 that they learn options that cause the system to reach states from which future subtasks are difficult
 250 to perform—e.g., in the rooms environment, the agent sometimes reaches the left half of the exits
 251 from where it is difficult to reach the right exit in the subsequent room. Although MADDPG uses the
 252 same reduction to two-player games as ROSAC, it ignores all the structure in the game and treats it as
 253 a generic Markov game. As a result, it learns a separate NN policy for each player which leads to the
 254 issue of unstable training, primarily due to the non-stationary nature of the environment observed
 255 by either agent. As shown in the plots, this leads to poor performance when applied to the problem
 256 of learning robust options.

257 6 Conclusions

258 We have proposed a framework for training robust options which can be used to perform arbitrary
 259 sequences of subtasks. In our framework, we first reduce the problem to a two-agent zero-sum
 260 stagewise Markov game which has a unique structure. We utilized this structure to design two al-
 261 gorithms, namely ROSAC and AROSAC, and demonstrated that they outperform existing approaches
 262 for training options with respect to multi-task performance. One potential limitation of our approach
 263 is that the set of subtasks is fixed and has to be provided by the user. An interesting direction for
 264 future work is to address this limitation by combining our approach with option discovery methods.

265 **Societal impacts.** Our work seeks to improve reinforcement learning for complex long-horizon
 266 tasks. Any progress in this direction would enable robotics applications both with positive impact—
 267 e.g., flexible and general-purpose manufacturing robotics, robots for achieving agricultural tasks,
 268 and robots that can be used to perform household chores—and with negative or controversial
 269 impact—e.g., military applications. These issues are inherent in all work seeking to improve the
 270 abilities of robots.

271 **References**

- 272 [1] David Abel, Nathan Umbanhowar, Khimya Khetarpal, Dilip Arumugam, Doina Precup, and
273 Michael L. Littman. Value preserving state-action abstractions. In *Proceedings of the Interna-*
274 *tional Conference on Artificial Intelligence and Statistics*, 2020.
- 275 [2] Natalia Akchurina. Multi-agent reinforcement learning algorithm with variable optimistic-
276 pessimistic criterion. In *ECAI*, volume 178, pages 433–437, 2008.
- 277 [3] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur
278 Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s
279 cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- 280 [4] Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with
281 policy sketches. In *International Conference on Machine Learning*, pages 166–175. PMLR,
282 2017.
- 283 [5] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Thirty-First*
284 *AAAI Conference on Artificial Intelligence*, 2017.
- 285 [6] Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In
286 *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- 287 [7] D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, Belmont,
288 MA, 1996.
- 289 [8] Christian Daniel, Herke Van Hoof, Jan Peters, and Gerhard Neumann. Probabilistic inference
290 for determining options in reinforcement learning. *Machine Learning*, 104(2):337–357, 2016.
- 291 [9] Thomas G Dietterich. State abstraction in maxq hierarchical reinforcement learning. In *Ad-*
292 *vances in Neural Information Processing Systems*, pages 994–1000, 2000.
- 293 [10] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you
294 need: Learning skills without a reward function. In *ICLR*, 2018.
- 295 [11] F1/10 Autonomous Racing Competition. <http://f1tenth.org>.
- 296 [12] Chelsea Finn, Tianhe Yu, Justin Fu, Pieter Abbeel, and Sergey Levine. Generalizing skills with
297 semi-supervised reinforcement learning. In *ICLR*, 2017.
- 298 [13] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in
299 actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596,
300 2018.
- 301 [14] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-
302 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*
303 *conference on machine learning*, pages 1861–1870. PMLR, 2018.
- 304 [15] Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Jour-*
305 *nal of machine learning research*, 4(Nov):1039–1069, 2003.
- 306 [16] Junling Hu, Michael P Wellman, et al. Multiagent reinforcement learning: theoretical frame-
307 work and an algorithm. In *ICML*, volume 98, pages 242–250. Citeseer, 1998.
- 308 [17] Radoslav Ivanov, Kishor Jothimurugan, Steve Hsu, Shaan Vaidya, Rajeev Alur, and Osbert
309 Bastani. Compositional learning and verification of neural network controllers. *ACM Trans-*
310 *actions on Embedded Computing Systems (TECS)*, 20(5s):1–26, 2021.
- 311 [18] Kishor Jothimurugan, Osbert Bastani, and Rajeev Alur. Abstract value iteration for hierarchi-
312 cal reinforcement learning. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings*
313 *of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of
314 *Proceedings of Machine Learning Research*, pages 1162–1170. PMLR, 13–15 Apr 2021.
- 315 [19] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European con-*
316 *ference on machine learning*, pages 282–293. Springer, 2006.

- 317 [20] George Konidaris and Andrew G Barto. Skill discovery in continuous reinforcement learning
318 domains using skill chaining. In *Advances in neural information processing systems*, pages
319 1015–1023, 2009.
- 320 [21] George Konidaris, Leslie Kaelbling, and Tomas Lozano-Perez. Constructing symbolic rep-
321 resentations for high-level planning. In *Twenty-Eighth AAAI Conference on Artificial Intelli-*
322 *gence*, 2014.
- 323 [22] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saedi, and Josh Tenenbaum. Hierarchical
324 deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In
325 *Advances in neural information processing systems*, pages 3675–3683, 2016.
- 326 [23] Yen-Ling Kuo, Boris Katz, and Andrei Barbu. Encoding formulas as deep networks: Re-
327 inforcement learning for zero-shot execution of ltl formulas. *2020 IEEE/RSJ International*
328 *Conference on Intelligent Robots and Systems (IROS)*, pages 5604–5610, 2020.
- 329 [24] Youngwoon Lee, Joseph J Lim, Anima Anandkumar, and Yuke Zhu. Adversarial skill chain-
330 ing for long-horizon robot manipulation via terminal state regularization. *arXiv preprint*
331 *arXiv:2111.07999*, 2021.
- 332 [25] Youngwoon Lee, Shao-Hua Sun, Sriram Somasundaram, Edward S Hu, and Joseph J Lim.
333 Composing complex skills by learning transition policies. In *International Conference on*
334 *Learning Representations*, 2018.
- 335 [26] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval
336 Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning.
337 In *ICLR*, 2016.
- 338 [27] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In
339 *Machine learning proceedings 1994*, pages 157–163. Elsevier, 1994.
- 340 [28] Michael L Littman. Friend-or-foe q-learning in general-sum games. In *ICML*, volume 1, pages
341 322–328, 2001.
- 342 [29] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent
343 actor-critic for mixed cooperative-competitive environments. In *Proceedings of the 31st Inter-*
344 *national Conference on Neural Information Processing Systems*, pages 6382–6393, 2017.
- 345 [30] Marios C Machado, Marc G Bellemare, and Michael Bowling. A laplacian framework for op-
346 tion discovery in reinforcement learning. In *Proceedings of the 34th International Conference*
347 *on Machine Learning-Volume 70*, pages 2295–2304. JMLR. org, 2017.
- 348 [31] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-optimal representation
349 learning for hierarchical reinforcement learning. In *ICLR*, 2019.
- 350 [32] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical
351 reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 3303–
352 3313, 2018.
- 353 [33] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generaliza-
354 tion with multi-task deep reinforcement learning. In *International Conference on Machine*
355 *Learning*, pages 2661–2670. PMLR, 2017.
- 356 [34] Stephen David Patek. *Stochastic and shortest path games: theory and algorithms*. PhD thesis,
357 Massachusetts Institute of Technology, 1997.
- 358 [35] Julien Perolat, Florian Strub, Bilal Piot, and Olivier Pietquin. Learning Nash Equilibrium
359 for General-Sum Markov Games from Batch Data. In *Proceedings of the 20th International*
360 *Conference on Artificial Intelligence and Statistics*, 2017.
- 361 [36] HL Prasad, Prashanth LA, and Shalabh Bhatnagar. Two-timescale algorithms for learning
362 nash equilibria in general-sum stochastic games. In *Proceedings of the 2015 International*
363 *Conference on Autonomous Agents and Multiagent Systems*, pages 1371–1379, 2015.

- 364 [37] Doina Precup, Richard S Sutton, and Satinder Singh. Theoretical results on reinforcement
365 learning with temporally abstract options. In *European conference on machine learning*, pages
366 382–393. Springer, 1998.
- 367 [38] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and
368 structured prediction to no-regret online learning. In *Proceedings of the fourteenth interna-*
369 *tional conference on artificial intelligence and statistics*, pages 627–635, 2011.
- 370 [39] Sungryull Sohn, Junhyuk Oh, and Honglak Lee. Hierarchical reinforcement learning for zero-
371 shot generalization with subtask dependencies. *Advances in Neural Information Processing*
372 *Systems*, 31, 2018.
- 373 [40] Martin Stolle and Doina Precup. Learning options in reinforcement learning. In *International*
374 *Symposium on abstraction, reformulation, and approximation*, pages 212–223. Springer, 2002.
- 375 [41] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A frame-
376 work for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–
377 211, 1999.
- 378 [42] Georgios Theodorou and Leslie P Kaelbling. Approximate planning in pomdps with macro-
379 actions. In *Advances in Neural Information Processing Systems*, pages 775–782, 2004.
- 380 [43] Saket Tiwari and Philip S Thomas. Natural option critic. In *Proceedings of the AAAI Confer-*
381 *ence on Artificial Intelligence*, volume 33, pages 5175–5182, 2019.
- 382 [44] Pashootan Vaezipoor, Andrew C Li, Rodrigo A Toro Icarte, and Sheila A. Mcilraith. Ltl2action:
383 Generalizing ltl instructions for multi-task rl. In Marina Meila and Tong Zhang, editors, *Pro-*
384 *ceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceed-*
385 *ings of Machine Learning Research*, pages 10497–10508. PMLR, 18–24 Jul 2021.
- 386 [45] Chen-Yu Wei, Yi-Te Hong, and Chi-Jen Lu. Online reinforcement learning in stochastic games.
387 In *Proceedings of the 31st International Conference on Neural Information Processing Sys-*
388 *tems*, pages 4994–5004, 2017.

389 Checklist

- 390 1. For all authors...
- 391 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
392 contributions and scope? [Yes]
- 393 (b) Did you describe the limitations of your work? [Yes]
- 394 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 395 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
396 them? [Yes]
- 397 2. If you are including theoretical results...
- 398 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 399 (b) Did you include complete proofs of all theoretical results? [Yes] In the supplement.
- 400 3. If you ran experiments...
- 401 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
402 mental results (either in the supplemental material or as a URL)? [Yes]
- 403 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
404 were chosen)? [Yes]
- 405 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
406 ments multiple times)? [Yes]
- 407 (d) Did you include the total amount of compute and the type of resources used (e.g., type
408 of GPUs, internal cluster, or cloud provider)? [Yes]
- 409 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 410 (a) If your work uses existing assets, did you cite the creators? [Yes]

- 411 (b) Did you mention the license of the assets? [N/A]
412 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
413
414 (d) Did you discuss whether and how consent was obtained from people whose data
415 you're using/curating? [N/A]
416 (e) Did you discuss whether the data you are using/curating contains personally identifi-
417 able information or offensive content? [N/A]
418 5. If you used crowdsourcing or conducted research with human subjects...
419 (a) Did you include the full text of instructions given to participants and screenshots, if
420 applicable? [N/A]
421 (b) Did you describe any potential participant risks, with links to Institutional Review
422 Board (IRB) approvals, if applicable? [N/A]
423 (c) Did you include the estimated hourly wage paid to participants and the total amount
424 spent on participant compensation? [N/A]