
Average-Reward Learning and Planning with Options

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We extend the options framework for temporal abstraction in reinforcement learning
2 from discounted Markov decision processes (MDPs) to average-reward MDPs.
3 Our contributions include general convergent off-policy inter-option learning algo-
4 rithms, intra-option algorithms for learning values and models, as well as sample-
5 based planning variants of our learning algorithms. Our algorithms and conver-
6 gence proofs extend those recently developed by Wan, Naik, and Sutton. We
7 also extend the notion of option-interrupting behaviour from the discounted to the
8 average-reward formulation. We show the efficacy of the proposed algorithms with
9 experiments on a continuing version of the Four-Room domain.

10 1 Introduction

11 Reinforcement learning (RL) is a formalism of trial-and-error learning in which an agent interacts
12 with an environment to learn a behavioral strategy that maximizes a notion of reward. In many
13 problems of interest, a learning agent may need to predict the consequences of its actions over
14 multiple levels of temporal abstraction. The *options* framework provides a way for defining courses
15 of actions over extended time scales, and for learning, planning, and representing knowledge with
16 them (Sutton, Precup, & Singh 1999, Sutton & Barto 2018). The options framework was originally
17 proposed within the *discounted* formulation of RL in which the agent tries to maximize the expected
18 discounted return from each state. We extend the options framework from the discounted formulation
19 to the *average-reward* formulation in which the goal is to find a policy that maximizes the rate of
20 reward.

21 Given a Markov decision process (MDP) and a fixed set of options, learning and planning algorithms
22 can be divided into two classes. The first class consists of *inter*-option algorithms, which enable an
23 agent to learn or plan with options instead of primitive actions. Given an option, the learning and
24 planning updates for this option in these algorithms occur only *after* the option’s actual or simulated
25 execution. Algorithms in this class are also called semi-MDP (SMDP) algorithms because given
26 an MDP, the decision process that selects among a set of options, executing each to termination,
27 is an SMDP (Sutton et al., 1999). The second class consists of algorithms in which learning or
28 planning updates occur after each state-action transition *within* options’ execution — these are called
29 *intra*-option algorithms. From a single state-action transition, these algorithms can learn or plan to
30 improve the values or policies for *all* options that may generate that transition, and are therefore
31 potentially more efficient than SMDP algorithms.

32 Several inter-option (SMDP) learning algorithms have been proposed for the average-reward formu-
33 lation (see, e.g., Das et al. 1999, Gosavi 2004, Vien & Chung 2008). To the best of our knowledge,
34 Gosavi’s (2004) algorithm is the only proven-convergent *off-policy* inter-option learning algorithm.
35 However, its convergence proof requires the underlying SMDP to have a special state that is recurrent
36 under all stationary policies. We extend Wan, Naik, and Sutton’s (2021) Differential Q-learning,
37 an off-policy control learning algorithm, to *inter-option Differential Q-learning* and show that it
38 converges without requiring a special state. For planning, we propose *inter-option Differential*

39 *Q-planning*, which is the first convergent *incremental* (sampled-based) planning algorithm. The
 40 existing proven-convergent inter-option planning algorithms (e.g., Schweitzer 1971, Puterman 1994,
 41 Li & Cao 2010) are not incremental because they perform a full sweep over states for each planning
 42 step.

43 Additionally, the literature seems to lack intra-option learning and planning algorithms within
 44 the average-reward formulation for both values and models. We fill this gap by proposing such
 45 algorithms in the average-reward formulation and provide their convergence results. These algorithms
 46 are stochastic approximation algorithms solving the average-reward intra-option value and model
 47 equations, which are also introduced in this paper for the first time.

48 Sutton et al. (1999) also introduced an algorithm to improve an agent’s behavior given estimated
 49 option values. Instead of letting an option execute to termination, this algorithm involves potentially
 50 interrupting an option’s execution to check if starting a new option might yield a better expected
 51 outcome. If so, then the currently-executing option is terminated, and the new option is executed. Our
 52 final contribution involves extending this notion of an *interruption* algorithm from the discounted to
 53 the average-reward formulation.

54 2 Problem Setting

55 We formalize an agent’s interaction with its environment by a finite Markov decision process (MDP)
 56 \mathcal{M} and a finite set of options \mathcal{O} . The MDP is defined by the tuple $\mathcal{M} \doteq (\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$, where \mathcal{S} is
 57 a set of states, \mathcal{A} is a set of actions, \mathcal{R} is a set of rewards, and $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the
 58 dynamics of the environment. Each option o in \mathcal{O} consists two components: the *option’s policy*
 59 $\pi^o : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, and a probability distribution of the *option’s termination* $\beta^o : \mathcal{S} \rightarrow [0, 1]$. For
 60 simplicity, for any $s \in \mathcal{S}, o \in \mathcal{O}$, we use $\pi(a|s, o)$ to denote $\pi^o(a, s)$ and $\beta(s, o)$ to denote $\beta^o(s)$.
 61 Sutton et al.’s (1999) options additionally have an *initiation* set that consists of the states at which
 62 the option can be initiated. To simplify the presentation in this paper, we allow all options to be
 63 initiated in all states of the state space; the algorithms and theoretical results can be easily extended
 64 to incorporate initiation from specific states.

65 An option o executes as follows. First, the agent observes a state S_t and chooses an action A_t
 66 according to the option’s policy $\pi(\cdot|S_t, o)$. The agent then observes the next state S_{t+1} and reward
 67 R_{t+1} according to p . The option either terminates at S_{t+1} with probability $\beta(S_{t+1}, o)$, or continues
 68 with action A_{t+1} chosen according to $\pi(\cdot|S_{t+1}, o)$. It then possibly terminates in S_{t+2} according
 69 to $\beta(S_{t+2}, o)$, and so on. The behavior of the agent is determined by a policy that chooses options,
 70 which we denote by $\mu_b : \mathcal{S} \times \mathcal{O} \mapsto [0, 1]$. In state S_t , the agent selects an option $O_t \in \mathcal{O}$ according
 71 to probability distribution $\mu_b(\cdot|S_t)$. The option policy starts executing at S_t and terminates S_{t+K} ,
 72 where K is a random variable denoting the number of time steps the option executed. At S_{t+K} , a
 73 new option is chosen according to $\mu_b(\cdot|S_{t+K})$, and so on. The agent-environment interactions go on
 74 forever without any resets. Note that actions are a special case of options—every action a is an option
 75 o that terminates after exactly one step ($\beta(s, o) = 1, \forall s$) and whose policy is to pick a in every state
 76 ($\pi(a|s, o) = 1, \forall s$).

77 Let T_n denote the time step when the $n - 1^{\text{th}}$ option terminates and the n^{th} option is chosen. Denote
 78 the n^{th} option by $\hat{O}_n \doteq O_{T_n}$, its starting state by $\hat{S}_n \doteq S_{T_n}$, the cumulative reward during its
 79 execution by $\hat{R}_n \doteq \sum_{t=T_n}^{T_{n+1}-1} R_t$, the state it terminates in by $\hat{S}_{n+1} \doteq S_{T_{n+1}}$, and its length by
 80 $\hat{L}_n \doteq T_{n+1} - T_n$. Note that every option’s length is a random variable taking values among positive
 81 integers. Then the option’s transition probability can be defined as $\hat{p}(s', r, l | s, o) \doteq \Pr(\hat{S}_{n+1} =$
 82 $s', \hat{R}_n = r, \hat{L}_n = l | \hat{S}_n = s, \hat{O}_n = o)$. Throughout the paper, we assume that the expected execution
 83 time of any option starting from any state is finite.

84 Given an MDP \mathcal{M} and the set of options \mathcal{O} , $\hat{\mathcal{M}} = (\mathcal{S}, \mathcal{O}, \hat{\mathcal{L}}, \hat{\mathcal{R}}, \hat{p})$ is a(an) SMDP, where $\hat{\mathcal{L}}$ is the set
 85 of all possible lengths of options and $\hat{\mathcal{R}}$ is the set of all possible options’ cumulative rewards. For this
 86 SMDP, the *reward rate* of a policy of interest, μ , given a starting state s and option o can be defined
 87 as $r^C(\mu)(s, o) \doteq \lim_{t \rightarrow \infty} \mathbb{E}_\mu[\sum_{i=1}^t R_i | S_0 = s, O_0 = o]/t$. Alternatively, at the level of option
 88 transitions, $r(\mu)(s, o) \doteq \lim_{n \rightarrow \infty} \mathbb{E}_\mu[\sum_{i=0}^n \hat{R}_i | \hat{S}_0 = s, \hat{O}_0 = o]/\mathbb{E}_\mu[\sum_{i=0}^n \hat{L}_i | \hat{S}_0 = s, \hat{O}_0 = o]$.
 89 It can be shown that the above two limits exist and are equivalent (Puterman’s (1994) propositions
 90 11.4.1 and 11.4.7) under the following assumption:

91 **Assumption 1.** Consider an MDP $(\mathcal{S}, \mathcal{O}, \hat{\mathcal{R}}, p')$, where $p'(s', r | s, o) \doteq \sum_l \hat{p}(s', r, l | s, o)$ for all
 92 s', r, s, o . The Markov chain induced by any stationary policy in this MDP is recurrent.

93 Under Assumption 1, the reward rate does not depend on the start state-option pair, and hence we
 94 denote it by just $r(\mu)$. The optimal reward rate can then be defined as $r_* \doteq \sup_{\mu \in \Pi} r(\mu)$, where Π
 95 denotes the set of all policies. The differential option-value function for a policy μ is defined for all
 96 $s \in \mathcal{S}, o \in \mathcal{O}$ as $q_\mu(s, o) \doteq \mathbb{E}_\mu[R_{t+1} - r(\mu) + R_{t+2} - r(\mu) + \dots | S_t = s, O_t = o]$. The *evaluation*
 97 and *optimality* equations for SMDPs are, as given by Puterman (1994):

$$q(s, o) = \sum_{s', r, l} \hat{p}(s', r, l | s, o) (r - \bar{r} \cdot l + \sum_{o'} \mu(o' | s') q(s', o')), \quad (1)$$

$$q(s, o) = \sum_{s', r, l} \hat{p}(s', r, l | s, o) (r - \bar{r} \cdot l + \max_{o'} q(s', o')). \quad (2)$$

98 Just like the average-reward MDP Bellman equations, the SMDP Bellman equations have a unique
 99 solution for $\bar{r} - r(\mu)$ for evaluation and r_* for control — and a unique solution for q only up to a
 100 constant. Given an MDP and a set of options, the goal of the *prediction* problem is, for a given policy
 101 μ , to find the reward rate $r(\mu)$ and the differential value function (possibly with some constant offset).
 102 The goal of the *control* problem is to find a policy that achieves the optimal reward rate r_* .

103 3 Inter-Option Learning and Planning Algorithms

104 In this section, we present our inter-option learning and planning, prediction and control algorithms,
 105 which extend Wan et al.'s (2021) differential learning and planning algorithms for average-reward
 106 MDPs from actions to options. We begin with the control learning algorithm and then move on to the
 107 prediction and planning algorithms.

108 Consider Wan et al.'s (2021) control learning algorithm:

$$Q_{t+1}(S_t, A_t) \doteq Q_t(S_t, A_t) + \alpha_t \delta_t, \quad \bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \delta_t,$$

109 where Q is a vector of size $|\mathcal{S} \times \mathcal{A}|$ that approximates a solution of q in the Bellman equation
 110 for MDPs, \bar{R} is a scalar estimate of the optimal reward rate, α_t is a step-size sequence, η is a
 111 positive constant, and δ_t is the temporal-difference (TD) error: $\delta_t \doteq R_t - \bar{R}_t + \max_a Q_t(S_{t+1}, a) -$
 112 $Q_t(S_t, A_t)$. The most straightforward inter-option extension of Differential Q-learning is:

$$Q_{n+1}(\hat{S}_n, \hat{O}_n) \doteq Q_n(\hat{S}_n, \hat{O}_n) + \alpha_n \delta_n, \quad (3)$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta \alpha_n \delta_n, \quad (4)$$

113 where Q is a vector of size $|\mathcal{S} \times \mathcal{O}|$ that approximates a solution of q in (2), \bar{R} is a scalar estimate of
 114 r_* , α_n is a step-size sequence, and δ_n is the TD error:

$$\delta_n \doteq \hat{R}_n - \hat{L}_n \bar{R}_n + \max_o Q_n(\hat{S}_{n+1}, o) - Q_n(\hat{S}_n, \hat{O}_n). \quad (5)$$

115 Such an algorithm is prone to instability because the option length \hat{L}_n can be quite large, and any
 116 error in the reward-rate estimate \bar{R}_n gets multiplied with the potentially-large option length. Using
 117 small step sizes might make the updates relatively stable, but at the cost of slowing down learning for
 118 options of shorter lengths. This could make the choice of step size quite critical, especially when the
 119 range of the options' lengths is large and unknown. Alternatively, inspired by Schweitzer (1971), we
 120 propose scaling the updates by the estimated length of the option being executed:

$$Q_{n+1}(\hat{S}_n, \hat{O}_n) \doteq Q_n(\hat{S}_n, \hat{O}_n) + \alpha_n \delta_n / L_n(\hat{S}_n, \hat{O}_n), \quad (6)$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta \alpha_n \delta_n / L_n(\hat{S}_n, \hat{O}_n), \quad (7)$$

121 where α_n is a step-size sequence, $L_n(\cdot, \cdot)$ comes from an additional vector of estimates $L : \mathcal{S} \times \mathcal{O} \rightarrow$
 122 \mathbb{R} that approximates expected lengths of state-option pairs, updated from experience by:

$$L_{n+1}(\hat{S}_n, \hat{O}_n) \doteq L_n(\hat{S}_n, \hat{O}_n) + \beta_n (\hat{L}_n - L_n(\hat{S}_n, \hat{O}_n)), \quad (8)$$

123 where β_n is an another step-size sequence. The TD-error δ_n in (6) and (7) is as in (5) with $L_n(\hat{S}_n, \hat{O}_n)$
 124 instead of \hat{L}_n . These equations make up our *inter-option Differential Q-learning* algorithm.

Similarly, our prediction learning algorithm, called *inter-option Differential Q-evaluation*, also has update rules (6–8) but with the TD error now equal to:

$$\delta_n \doteq \hat{R}_n - L_n(\hat{S}_n, \hat{O}_n) \bar{R}_n + \sum_o \mu(o|\hat{S}_{n+1}) Q_n(\hat{S}_{n+1}, o) - Q_n(\hat{S}_n, \hat{O}_n). \quad (9)$$

Theorem 1 (Convergence of intra-option algorithms, informal). *If Assumption 1 holds, step sizes are decreased appropriately, all state-option pairs are visited for an infinite number of times, and the relative visitation frequency between any two pairs is finite:*

1. *inter-option Differential Q-learning (5–8) converges almost surely, \bar{R}_n to r_* and Q_n to a solution of (2), and $r(\mu_n)$ to r_* , where μ_n is a greedy policy w.r.t. Q_n ,*
2. *inter-option Differential Q-evaluation (6–9) converges almost surely, \bar{R}_n to $r(\mu)$ and Q_n to a solution of (1).*

The convergence proofs for the inter-option (as well as the subsequent intra-option) algorithms are based on a result that generalizes Wan et al.’s (2021) and Abounadi et al.’s (2001) proof techniques from primitive actions to options. We present this result in Appendix A.1 and the formal theorem statements as well as proofs in Appendix A.2.

Remark: The scaling factor $L_n(\hat{S}_n, \hat{O}_n)$ used in the algorithm is the expected option length instead of the sampled option length. Scaling the updates by the expected option lengths guarantees that fixed points of the updates are the same as those of (3–4), which are the solutions of (2). This is not guaranteed to be true when using the sampled option length. We discuss this in more detail in Appendix C.1.

The inter-option planning algorithms for prediction and control are similar to the learning algorithms except that they use simulated experience generated by a (given or learned) model instead of real experience. In addition, they only have two update rules, (6) and (7), and not (8) because the model provides the expected length of a given option from a given state (see Section 5 for a complete specification of option models). The planning algorithms and their convergence results are presented in Appendix A.2.

Empirical Evaluation. We tested our inter-option Differential Q-learning with Gosavi’s (2004) algorithm as a baseline in a variant of Sutton et al.’s (1999) Four-Room domain (shown in Figure 1). The agent starts in the yellow cell. The goal states are indicated by green cells. Every experiment in this paper uses only one of the green cells as a goal state; the other two are considered as empty cells. There are four primitive actions of moving up, down, left, right. The agent receives a reward of +1 when it moves into the goal cell, 0 otherwise.

In addition to the four primitive actions, the agent has eight options that take it from a given room to the hallway adjoining the room. The arrows in Figure 1 illustrate the policy of one of the eight options. For this option, the policy in the empty cells (not marked with arrows) is to uniformly-randomly pick among the four primitive actions. The termination probability is 0 for all the cells with arrows and 1 for the empty cells. The other seven options are defined in a similar way. Denote the set of primitive actions as \mathcal{A} and the set of hallway options as \mathcal{H} . Including the primitive actions, the agent has 12 options in total.

In the first experiment, we tested inter-option Differential Q-learning with three different sets of options, $\mathcal{O} = \mathcal{A}$, $\mathcal{O} = \mathcal{H}$ and $\mathcal{O} = \mathcal{A} + \mathcal{H}$. The task was to reach the green cell G1, the shortest path to which from the starting state is 16 steps. Hence the best possible reward rate for this task is $1/16 = 0.0625$. The agent used an ϵ -greedy policy with $\epsilon = 0.1$. For each of the two step-sizes α_n and β_n , we tested five choices: 2^{-x} , $x \in \{1, 3, 5, 7, 9\}$. In addition, we tested five choices of $\eta : 10^{-x}$, $x \in \{0, 1, 2, 3, 4\}$. Q and \bar{R} were initialized to 0; L was initialized to 1. Each parameter setting was run for 200,000 steps and repeated 30 times. The left subfigure of Figure 2 shows a typical learning curve for each of the three sets of options, with $\alpha = 2^{-3}$, $\beta = 2^{-1}$, and $\eta = 10^{-1}$. The parameter study for $\mathcal{O} = \mathcal{A} + \mathcal{H}$ w.r.t. α and η , with $\beta = 2^{-1}$, is presented in the right subfigure of Figure 2. The metric is the average reward obtained over the entire training period. Complete parameter studies for $\mathcal{O} = \mathcal{H} + \mathcal{A}$, $\mathcal{O} = \mathcal{H}$ and $\mathcal{O} = \mathcal{A}$ are presented in Appendix B.1.

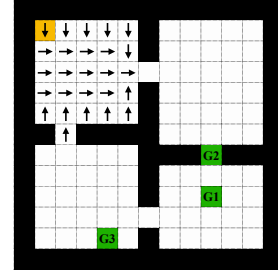


Figure 1: A continuing variant of the Four-Room domain where the task is to repetitively go from the yellow start state to one of the three green goal states. Also shown is an option policy to go to the upper hallway cell; more details in-text.

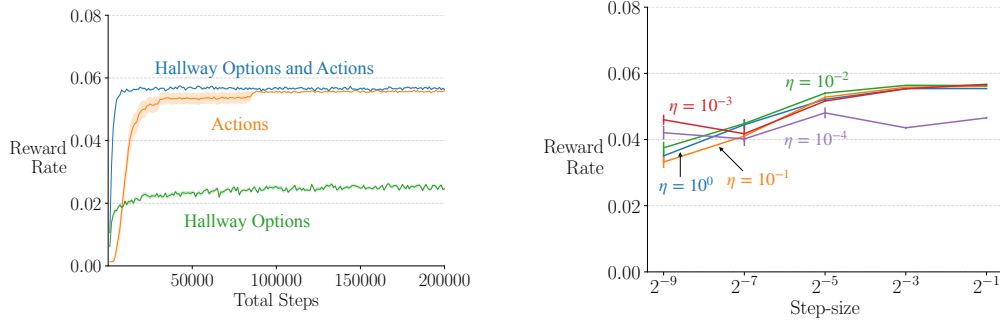


Figure 2: Plots showing some learning curves and the parameter study of inter-option Differential Q-learning on the continuing Four-Room domain when the goal was to go to G1. *Left*: A point on the solid line denotes reward rate over the last 1000 time steps and the shaded region indicates one standard error. The behavior using the three different sets of options was as expected. *Right*: Sensitivity of performance to α and η when using $\mathcal{O} = \mathcal{A} + \mathcal{H}$ and $\beta = 2^{-1}$. The x-axis denotes step size α ; the y-axis denotes the rate of the rewards averaged over all 200,000 steps of training, reflecting the rate of learning. The error bars denote one standard error. The algorithm’s rate of learning varied little over a broad range of its parameters α and η .

178 The learning curves in the left panel of Figure 2 show that the agent achieved a relatively stable
 179 reward rate after 100,000 steps in all three cases. Using just primitive actions, the learning curve rises
 180 the slowest, indicating that hallway options indeed help the agent reach the goal faster. But solely
 181 using the hallway options is not very useful in the long run as the goal G1 is not a hallway state. Note
 182 that because of the ϵ -greedy behavior policy, the learning curves do not reach the optimal reward rate
 183 of 0.0625. These observations mirror those by Sutton et al. (1999) in the discounted formulation.

184 The sensitivity curves of inter-option Differential Q-learning (right panel of Figure 2) indicate that,
 185 in this Four-Room domain, the algorithm was not sensitive to parameters η , performed well for a
 186 wide range of step sizes α , and showed low variance across different runs. We also found that the
 187 algorithm was not sensitive to β either; this and the additional parameter studies involving the two
 188 other option sets are presented in Appendix B.1.

189 We also tested Gosavi’s (2004) algorithm as a baseline.
 190 Recall it is the only proven-convergent SMDP off-policy
 191 control learning algorithm prior to our work. The algo-
 192 rithm estimates the reward rate by tracking the cumulative
 193 reward \bar{C} obtained by the options and dividing it by the
 194 another estimate \bar{T} the tracks the length of the options.
 195 If the n^{th} option executed is a greedy choice, then these
 196 estimates are updated using:

$$\begin{aligned}\bar{C}_{n+1} &\doteq \bar{C}_n + \beta_n(\hat{R}_n - C_n), \\ \bar{T}_{n+1} &\doteq \bar{T}_n + \beta_n(\hat{L}_n - T_n), \\ \bar{R}_{n+1} &\doteq \bar{C}_{n+1}/\bar{T}_{n+1}.\end{aligned}$$

197 When \hat{O}_n is not greedy, $\bar{R}_{n+1} = \bar{R}_n$. The option-value
 198 function is updated with (3) with δ_n as defined in (5). α_n
 199 and β_n are two step-size sequences. The sensitivity of
 200 this algorithm with $\mathcal{O} = \mathcal{A} + \mathcal{H}$ is shown in Figure 3.

201 Compared to inter-option Differential Q-learning, this baseline has one less parameter, but its
 202 performance was found to be more sensitive to the values of both its step-size parameters. In addition,
 203 the error bars were generally larger, suggesting that the variance across different runs was also higher.

204 To conclude, our experiments with the continuing Four-Room domain show that our inter-option
 205 Differential Q-learning indeed finds the optimal policy given a set of options, in accordance with
 206 Theorem 1. In addition, its performance seems to be robust to the choices of its parameters.

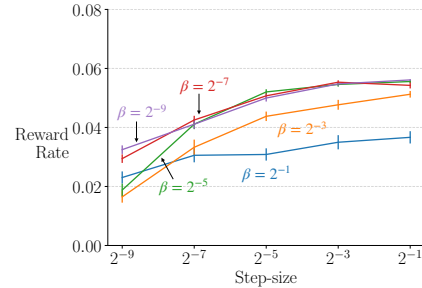


Figure 3: Parameter studies showing Gosavi’s (2004) algorithm’s rate of learning is relatively more sensitive to the choices of its two parameters compared to our inter-option Differential Q-learning. The experimental setting and the plot axes are the same as mentioned in Figure 2’s caption.

4 Intra-Option Value Learning and Planning Algorithms

In this section, we introduce intra-option value learning and planning algorithms. The objectives are same as that of inter-option value learning algorithms. As mentioned earlier, intra-option algorithms learn from every transition $S_t, A_t, R_{t+1}, S_{t+1}$ during the execution of a given option O_t . Moreover, intra-option algorithms also make updates for every option $o \in \mathcal{O}$, including ones that may potentially never be executed.

To develop our algorithms, we first establish the intra-option evaluation and optimality equations in the average-reward case. The general form of the intra-option Bellman equation is:

$$q(s, o) = \sum_a \pi(a | s, o) \sum_{s', r} p(s', r | s, a) (r - \bar{r} + u^q(s', o)) \quad (10)$$

where $q \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}|}$ and $\bar{r} \in \mathbb{R}$ are free variables. The optimality and evaluation equations use $u^q = u_*^q$ and $u^q = u_\mu^q$ respectively, defined $\forall s' \in \mathcal{S}, o \in \mathcal{O}$ as:

$$u^q(s', o) = u_*^q(s', o) \doteq (1 - \beta(s', o))q(s', o) + \beta(s', o) \max_{o'} q(s', o'), \quad (11)$$

$$u^q(s', o) = u_\mu^q(s', o) \doteq (1 - \beta(s', o))q(s', o) + \beta(s', o) \sum_{o'} \mu(o' | s') q(s', o'). \quad (12)$$

Intuitively, the u^q term accounts for the two possibilities of an option terminating or continuing in the next state. These equations generalize the average-reward Bellman equations given by Puterman (1994). The following theorem characterizes the solutions to the intra-option Bellman equations.

Theorem 2 (Solutions to intra-option Bellman equations). *If Assumption 1 holds, then:*

1. a) there exists a $\bar{r} = \mathbb{R}$ and a $q \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}|}$ for which (10) and (11) holds, b) the solution of \bar{r} is unique and is equal to r_* , the solutions of q form a set $\{q_* + ce \mid c \in \mathbb{R}\}$ where e is an all-one vector of size $|\mathcal{S}| \times |\mathcal{O}|$, c) a greedy policy w.r.t. a solution of q achieves the optimal reward rate r_* .
2. a) there exists a $\bar{r} \in \mathbb{R}$ and a $q \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}|}$ for which (10) and (12) holds, b) the solution of \bar{r} is unique and is equal to $r(\mu)$, the solutions of q form a set $\{q_\mu + ce \mid c \in \mathbb{R}\}$.

The proof extends that of Corollary 8.2.7, Theorem 8.4.3, and Theorem 8.4.4 by Puterman (1994), and is presented in Appendix A.3.

Our intra-option control and prediction algorithms are stochastic approximation algorithms solving the intra-option optimality and evaluation equations respectively. Both the algorithms maintain a vector of estimates $Q(s, o)$ and a scalar estimate \bar{R} , just like our inter-option algorithms (since intra-option algorithms make updates after every transition, they do not need to maintain an estimator for option lengths (L) like the inter-option algorithms). Our control algorithm, called *intra-option Differential Q-learning*, updates estimates Q and \bar{R} by:

$$Q_{t+1}(S_t, o) \doteq Q_t(S_t, o) + \alpha_t \rho_t(o) \delta_t(o), \quad \forall o \in \mathcal{O}, \quad (13)$$

$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \sum_{o \in \mathcal{O}} \rho_t(o) \delta_t(o), \quad (14)$$

where α_t is a step-size sequence, $\rho_t(o) \doteq \frac{\pi(A_t | S_t, o)}{\pi(A_t | S_t, O_t)}$ is the importance sampling ratio, and:

$$\delta_t(o) \doteq R_{t+1} - \bar{R}_t + u_*^{Q_t}(S_{t+1}, o) - Q_t(S_t, o). \quad (15)$$

Our prediction algorithm, called *intra-option Differential Q-evaluation*, also update Q and \bar{R} by (13) and (14) but with the TD error:

$$\delta_t(o) \doteq R_{t+1} - \bar{R}_t + u_\mu^{Q_t}(S_{t+1}, o) - Q_t(S_t, o). \quad (16)$$

Theorem 3 (Convergence of intra-option algorithms; informal). *Under the same assumptions as those of Theorem 1:*

1. *intra-option Differential Q-learning algorithm (13–15) converges almost surely, \bar{R}_t to r_* , Q_t to a solution of q in (10) and (11), and $r(\mu_t)$ to r_* , where μ_t is a greedy policy w.r.t. Q_t ,*

2. *intra-option Differential Q-evaluation algorithm (13,14,16) converges almost surely, \bar{R}_t to $r(\mu)$, Q_t to a solution of q in (10) and (12).*

Remark: The intra-option learning methods introduced in this section can be used with options having stochastic policies. This is possible with the use of the important sampling ratios as described above. Sutton et al.’s (1999) discounted intra-option learning methods were restricted to options having deterministic policies.

Again, the intra-option value planning algorithms are similar to the learning algorithms except that they use simulated experience generated by a given or learned model instead of real experience. The planning algorithms and their convergence results are presented in Appendix A.4.

Empirical Evaluation. We conducted another experiment in the Four-Room domain to show that intra-option Differential Q-learning can learn the values of hallway options \mathcal{H} using only primitive actions \mathcal{A} . As mentioned earlier, there are no intra-option average-reward baseline algorithms, so this is a proof-of-concept experiment.

The goal state for this experiment was G2, which can be reached using the option that leads to the lower hallway. The optimal reward rate in this case is $1/14 \approx 0.714$ with both $\mathcal{O} = \mathcal{H}$ and $\mathcal{O} = \mathcal{A}$. We applied intra-option Differential Q-learning using a behavior policy that chose the four primitive actions with equal probability in all states. Each parameter setting was run for 200,000 steps and repeated 30 times. For evaluation, we saved the learned option value function after every 1000 steps and computed the average reward of the corresponding greedy policy over 1000 steps.

Figure 4 shows the learning curve of this average reward across the 30 independent runs for parameters $\alpha = 0.125, \eta = 0.1$. The agent indeed succeeds in learning the option-value function corresponding to the hallway options using a behavior policy consisting only of primitive actions. The parameter study of intra-option Differential Q-learning is presented in Appendix B.2.

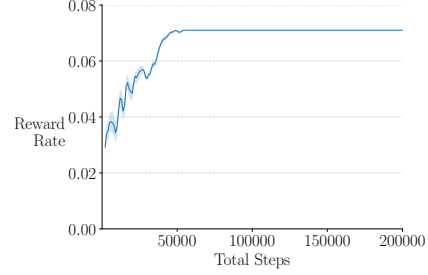


Figure 4: Learning curve showing that the greedy policy corresponding to the hallway options’ option-value function achieves the optimal reward rate on the continuing Four-Room domain. The value function was learned via intra-option Differential Q-learning using a behavior policy consisting only of primitive actions; the hallway options were never executed.

5 Intra-Option Model Learning and Planning Algorithms

In this section, we first describe option models within the average-reward formulation. We then introduce an algorithm to learn such models in an intra-option fashion. This option-model learning algorithm can be combined with the planning algorithms from the previous section to obtain a complete model-based average-reward options algorithm that learns option models and plans with them (we present this algorithm in Appendix C.2).

The average-reward option model is similar to the discounted options model but with key distinctions. Sutton et al.’s (1999) discounted option model has two parts: the dynamics part and the reward part. Given a state and an option, the dynamics part predicts the discounted occupancy of states upon termination, and the reward part predicts the expected (discounted) sum of rewards during the execution of the option. In the average-reward setting, apart from the dynamics and the reward parts, an option model has a third part—the *duration* part—that predicts the duration of execution of the option. In addition, the dynamics part predicts the state distribution upon termination without discounting and reward part predicts the undiscounted cumulative rewards during the execution of the option.

Formally, the dynamics part is $m^p(s'|s, o) \doteq \sum_{r,l} \hat{p}(s', r, l | s, o)$, the probability that option o terminates in state s' when starting from state s . The reward part is $m^r(s, o) \doteq \sum_{r,l} \hat{p}(s', r, l | s, o) r$, the expected cumulative reward during the execution of option o when starting from state s . Finally, the duration part is $m^l(s, o) \doteq \sum_{r,s'} \hat{p}(s', r, l | s, o) l$, the expected duration of option o when starting from state s .

We now present a set of recursive equations that are key to our model-learning algorithms. These equations extend the discounted Bellman equations for option models (Sutton et al. 1999) to the average-reward formulation.

$$\bar{m}^p(x | s, o) = \sum_a \pi(a | s, o) \sum_r p(s', r | s, a) (\beta(s', o) \mathbb{I}(x = s') + (1 - \beta(s', o)) \bar{m}^p(x | s', o)), \quad (17)$$

$$\bar{m}^r(s, o) = \sum_a \pi(a | s, o) \sum_{s', r} p(s', r | s, a) (r + (1 - \beta(s', o)) \bar{m}^r(s', o)), \quad (18)$$

$$\bar{m}^l(s, o) = \sum_a \pi(a | s, o) \sum_{s', r} p(s', r | s, a) (1 + (1 - \beta(s', o)) \bar{m}^l(s', o)). \quad (19)$$

The following theorem (see Appendix A.5 for the proof) shows that (m^p, m^r, m^l) is the unique solution of (17–19) and therefore the models can be obtained by solving these equations.

Theorem 4 (Solutions to Bellman equations for option models). *There exist unique $\bar{m}^p \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}| \times |\mathcal{S}|}$, $\bar{m}^r \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}|}$, and $\bar{m}^l \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}|}$ for which (17), (18), and (19) hold. Further, if $\bar{m}^p, \bar{m}^r, \bar{m}^l$ satisfy (17), (18), and (19), then $\bar{m}^p = m^p, \bar{m}^r = m^r, \bar{m}^l = m^l$.*

Our intra-option model-learning algorithm solves the above recursive equations using the following TD-like update rules for each option o :

$$\begin{aligned} M_{t+1}^p(x | St, o) &\doteq M_t^p(x | St, o) + \alpha_t \rho_t(o) \left(\beta(St_{t+1}, o) \mathbb{I}(St_{t+1} = x) \right. \\ &\quad \left. + (1 - \beta(St_{t+1}, o)) M_t^p(x | St_{t+1}, o) - M_t^p(x | St, o) \right), \quad \forall x \in \mathcal{S}, \end{aligned} \quad (20)$$

$$M_{t+1}^r(St, o) \doteq M_t^r(St, o) + \alpha_t \rho_t(o) \left(R_{t+1} + (1 - \beta(St_{t+1}, o)) M_t^r(St_{t+1}, o) - M_t^r(St, o) \right) \quad (21)$$

$$M_{t+1}^l(St, o) \doteq M_t^l(St, o) + \alpha_t \rho_t(o) \left(1 + (1 - \beta(St_{t+1}, o)) M_t^l(St_{t+1}, o) - M_t^l(St, o) \right) \quad (22)$$

where M^p is a $|\mathcal{S}| \times |\mathcal{O}| \times |\mathcal{S}|$ -sized vector of estimates, M^r and M^l are both $|\mathcal{S}| \times |\mathcal{O}|$ -sized vectors of estimates, and α_t is a sequence of step sizes. Standard stochastic approximation results can be applied to show the algorithm’s convergence (see Appendix A.6 for details).

Theorem 5 (Convergence of the intra-option model learning algorithm, informal). *Under Assumption 1, if the step sizes are set appropriately and all the state-option pairs are updated an infinite number of times, then the intra-option model-learning algorithm (20–22) converges almost surely, M_t^p to m^p , M_t^r to m^r , and M_t^l to m^l .*

Our intra-option model-learning algorithms (20–22) can be applied with simulated one-step transitions generated by a *given action model*, resulting to a planning algorithm that produces an *estimated option model*. The planning algorithm and its convergence result are presented in Appendix A.6.

6 Interruption to Improve Policy Over Options

In all the algorithms we considered so far, the policy over options would select an option, execute the option policy till termination, then select a new option. Sutton et al. (1999) showed that the policy over options can be improved by allowing the *interruption* of an option midway through its execution to start a seemingly better option. We now show that this interruption result applies for average-reward options as well (see Appendix A.7 for the proof).

Theorem 6 (Interruption). *For any MDP, any set of options \mathcal{O} , and any policy $\mu : \mathcal{S} \times \mathcal{O} \rightarrow [0, 1]$, define a new set of options, \mathcal{O}' , with a one-to-one mapping between the two option sets as follows: for every $o = (\pi, \beta) \in \mathcal{O}$, define a corresponding $o' = (\pi, \beta') \in \mathcal{O}'$ where $\beta' = \beta$, but for any state s in which $q_\mu(s, o) < v_\mu(s)$, $\beta'(s) = 1$. Let the interrupted policy μ' be such that for all $s \in \mathcal{S}$ and for all $o' \in \mathcal{O}'$, $\mu'(s, o') = \mu(s, o)$, where o is the option in \mathcal{O} corresponding to o' . Then:*

1. the new policy over options μ' is not worse than the old one μ , i.e., $r(\mu') \geq r(\mu)$,
2. if there exists a state $s \in \mathcal{S}$ from which there is a non-zero probability of encountering an interruption upon initiating μ' in s , then $r(\mu') > r(\mu)$.

In short, the above theorem shows that interruption produces a behavior that achieves a higher reward rate than without interruption. Note that interruption behavior is only applicable with intra-option algorithms; complete option transitions are needed in inter-option algorithms.

Empirical Evaluation. We tested the intra-option Differential Q-learning algorithm with and without interruption in the Four-Room domain. We set the goal as G3 and allowed the agent to choose and learn only from the set of all hallway options \mathcal{H} . With just hallway options, without interruption, the best strategy is to first move to the lower hallway and then try to reach the goal by following options that pick random actions in the states near the hallway and goal. With interruption, the agent can first move to the left hallway, then take the option that moves the agent to the lower hallway but terminate when other options have higher option-values. This termination is most likely to occur in the cell just above G3. The agent then needs a fewer number of steps in expectation to reach the goal.

Figure 5 shows learning curves using intra-option Differential Q-learning with and without interruptions on this problem. Each parameter setting was run for 400,000 steps and repeated 30 times. The learning curves shown correspond to $\alpha = 0.125$ and $\eta = 0.1$. As expected, the agent achieved a higher reward rate by using interruptions. The parameter study of the interruption algorithm along with the rest of the experimental details is presented in Appendix B.3.

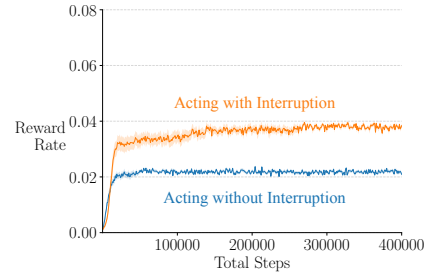


Figure 5: Learning curves showing that executing options with interruptions can achieve a higher reward rate than executing options till termination in the domain described in the adjoining text.

7 Conclusions, Limitations, and Future Work

In this paper, we extended learning and planning algorithms for the options framework — originally proposed by Sutton et al. (1999) for discounted-reward MDPs — to average-reward MDPs. The inter-option learning algorithm presented in this paper is more general than previous work in that its convergence proof does not require existence of any special states in the MDP. We also established intra-option Bellman equations in average-reward MDPs and used them to propose the first intra-option learning algorithms for average-reward MDPs. Finally, we extended the interruption algorithm and its related theory from the discounted to the average-reward setting. Our experiments on a continuing version of the classic Four-Room domain show the efficacy of the proposed algorithms. We believe that our contributions will enable widespread use of options in the average-reward setting.

The most immediate line of future work involves extending these ideas from the tabular case to the general case of function approximation, starting with linear function approximation. The idea of linear (discounted) options can be extended to the average-reward case, perhaps by building on the theory used by Zhang et al. (2021). Using the results developed in this paper, we also foresee extensions to more ideas from the discounted setting involving function approximation such as Bacon et al.’s (2017) option-critic architecture to the average-reward setting.

This paper assumes that a fixed set of options is provided and the agent then learns or plans using them. One of the most important challenges in the options framework is the *discovery* of options. We think the discovery problem is orthogonal to the problem formulation. Hence existing option-discovery algorithms developed for the discounted setting (e.g., by McGovern & Barto 2001, Menache et al. 2002, Şimşek & Barto 2004, Singh et al. 2004, Van Dijk & Polani 2011, Machado et al. 2017) can be easily extended to the average-reward setting. Relatively more work might be required in extending approaches that couple the problems of option discovery and learning (e.g., Gregor et al. 2016, Eysenbach et al. 2018, Achiam et al. 2018, Veeriah et al. 2021).

Another limitation of this paper is that it deals with learning and planning separately. We also need combined methods that learn models and plan with them; some ideas are discussed in Appendix C. Finally, we would like to get more empirical experience with the algorithms proposed in this paper, both in pedagogical tabular problems and challenging large-scale problems. Nevertheless, we believe this paper makes novel contributions that are significant for the use of temporal abstractions in average-reward reinforcement learning.

References

- Abounadi, J., Bertsekas, D., & Borkar, V. S. (2001). Learning Algorithms for Markov Decision Processes with Average Cost. *SIAM Journal on Control and Optimization*.
- Achiam, J., Edwards, H., Amodei, D., & Abbeel, P. (2018). Variational Option Discovery Algorithms. *ArXiv:1807.10299*.
- Almezel, S., Ansari, Q. H., & Khamsi, M. A. (2014). *Topics in Fixed Point Theory (Vol. 5)*. Springer.
- Bacon, P. L., Harb, J., & Precup, D. (2017). The Option-Critic Architecture. *AAAI Conference on Artificial Intelligence*.
- Borkar, V. S. (1998). Asynchronous Stochastic Approximations. *SIAM Journal on Control and Optimization*.
- Borkar, V. S. (2009). *Stochastic Approximation: A Dynamical Systems Viewpoint*. Springer.
- Borkar, V. S., & Soumyanatha, K. (1997). An Analog Scheme for Fixed Point Computation. I. Theory. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*.
- Brunskill, E., & Li, L. (2014). PAC-inspired Option Discovery in Lifelong Reinforcement Learning. *International Conference on Machine Learning*.
- Das, T. K., Gosavi, A., Mahadevan, S., & Marchallick, N. (1999). Solving Semi-Markov Decision Problems Using Average Reward Reinforcement Learning. *Management Science*.
- Eysenbach, B., Gupta, A., Ibarz, J., & Levine, S. (2018). Diversity is All You Need: Learning Skills without a Reward Function. *ArXiv:1802.06070*.
- Fruit, R., & Lazaric, A. (2017). Exploration-Exploitation in MDPs with Options. *Artificial Intelligence and Statistics*.
- Gosavi, A. (2004). Reinforcement Learning for Long-run Average Cost. *European Journal of Operational Research*.
- Gregor, K., Rezende, D. J., & Wierstra, D. (2016). Variational Intrinsic Control. *ArXiv:1611.07507*.
- Li, Y., & Cao, F. (2010). RVI Reinforcement Learning for Semi-Markov Decision Processes with Average Reward. *IEEE World Congress on Intelligent Control and Automation*.
- Machado, M. C., Bellemare, M. G., & Bowling, M. (2017). A Laplacian Framework for Option Discovery in Reinforcement Learning. *International Conference on Machine Learning*.
- McGovern, A., & Barto, A. G. (2001). Automatic Discovery of Subgoals in Reinforcement Learning using Diverse Density. *International Conference on Machine Learning*.
- Menache, I., Mannor, S., & Shimkin, N. (2002). Q-Cut - Dynamic Discovery of Sub-goals in Reinforcement Learning. *European Conference on Machine Learning*.
- Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons.
- Schweitzer, P. J. (1971). Iterative solution of the functional equations of undiscounted Markov renewal programming. *Journal of Mathematical Analysis and Applications*.
- Schweitzer, P. J., & Federgruen, A. (1978). The Functional Equations of Undiscounted Markov Renewal Programming. *Mathematics of Operations Research*.
- Şimşek, Ö., & Barto, A. G. (2004). Using Relative Novelty to Identify Useful Temporal Abstractions in Reinforcement Learning. *International Conference on Machine Learning*.
- Singh, S., Barto, A. G., & Chentanez, N. (2004). Intrinsically Motivated Reinforcement Learning. *Advances in Neural Information Processing Systems*.
- Sorg, J., & Singh, S. (2010). Linear Options. *International Conference on Autonomous Agents and Multiagent Systems*.
- Sutton, R. S., Precup, D., & Singh, S. (1999). Between MDPs and Semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning. *Artificial Intelligence*.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.

428 Tsitsiklis, J. N. (1994). Asynchronous stochastic approximation and Q-learning. *Machine Learning*.
 429 van Dijk, S. G. & Polani, D. (2011). Grounding Subgoals in Information Transitions. *IEEE*
 430 *Symposium on Adaptive Dynamic Programming and Reinforcement Learning*.
 431 Veeriah V., Zahavy T., Hessel M., Xu Z., Oh J., Kemaev I., van Hasselt H., Silver D., & Singh S.
 432 (2021). Discovery of Options via Meta-Learned Subgoals. *ArXiv:2102.06741*.
 433 Vien, N. A., & Chung, T. (2008). Policy Gradient Semi-Markov Decision Process. *IEEE International*
 434 *Conference on Tools with Artificial Intelligence*.
 435 Wan, Y., Naik, A., & Sutton, R. S. (2021). Learning and Planning in Average-Reward Markov
 436 Decision Processes. (To appear in) *International Conference on Machine Learning*.
 437 Yao, H., Szepesvári, C., Sutton, R. S., Modayil, J., & Bhatnagar, S. (2014). Universal Option Models.
 438 *Advances in Neural Information Processing Systems*.
 439 Zhang, S., Wan, Y., Sutton, R. S., & Whiteson, S. (2021). Average-Reward Off-Policy Policy
 440 Evaluation with Function Approximation. (To appear in) *International Conference on Machine*
 441 *Learning*.

442 Checklist

- 443 1. For all authors...
 - 444 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
 445 contributions and scope? [Yes]
 - 446 (b) Did you describe the limitations of your work? [Yes]
 - 447 (c) Did you discuss any potential negative societal impacts of your work? [No]
 - 448 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 449 them? [Yes]
- 450 2. If you are including theoretical results...
 - 451 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Formal
 452 statements of theoretical results can be found in Appendix A.
 - 453 (b) Did you include complete proofs of all theoretical results? [Yes] Formal proofs of
 454 theoretical results can be found in Appendix A.
- 455 3. If you ran experiments...
 - 456 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 457 mental results (either in the supplemental material or as a URL)? [No]
 - 458 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 459 were chosen)? [Yes]
 - 460 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 461 ments multiple times)? [Yes]
 - 462 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 463 of GPUs, internal cluster, or cloud provider)? [No]
- 464 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 465 (a) If your work uses existing assets, did you cite the creators? [N/A]
 - 466 (b) Did you mention the license of the assets? [N/A]
 - 467 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - 468
 - 469 (d) Did you discuss whether and how consent was obtained from people whose data you're
 470 using/curating? [N/A]
 - 471 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 472 information or offensive content? [N/A]
- 473 5. If you used crowdsourcing or conducted research with human subjects...
 - 474 (a) Did you include the full text of instructions given to participants and screenshots, if
 475 applicable? [N/A]

- 476 (b) Did you describe any potential participant risks, with links to Institutional Review
477 Board (IRB) approvals, if applicable? [N/A]
- 478 (c) Did you include the estimated hourly wage paid to participants and the total amount
479 spent on participant compensation? [N/A]

A Formal Theoretical Results and Proofs

In this section, we provide formal statements of the theorems presented in the main text of the paper and show their proofs. This section has several subsections. The first subsection introduces General RVI Q, which will be used in later subsections. The other six subsections correspond to six theorems presented in the main text.

A.1 General RVI Q

Wan et al. (2021) extended the family of RVI Q-learning algorithms (Abounadi, Bertsekas, and Borkar et al. 2001) to prove the convergence of their Differential Q-learning algorithm. Unlike RVI Q-learning, Differential Q-learning does not require a reference function. We further extend Wan et al.'s extended family of RVI Q-learning algorithms to a more general family of algorithms, called *General RVI Q*. We then prove convergence for this family of algorithms and show that inter-option algorithms and intra-option value learning algorithms are all members of this family.

We first need the following definitions:

1. a set-valued process $\{Y_n\}$ taking values in the set of nonempty subsets of \mathcal{I} with the interpretation: $Y_n = \{i : i^{\text{th}} \text{ component of } Q \text{ was updated at time } n\}$,
2. $\nu(n, i) \doteq \sum_{k=0}^n I\{i \in Y_k\}$, where I is the indicator function. Thus $\nu(n, i)$ = the number of times the i component was updated up to step n ,
3. i.i.d. random vectors R_n, G_n and F_n for all $n \geq 0$ satisfying $\mathbb{E}[R_n(i)] = r(i)$, where r is a fixed real vector, $\mathbb{E}[G_n(Q)(i)] = g(Q)(i)$ for any $Q \in \mathbb{R}^{|\mathcal{I}|}$ where $g : \mathcal{I} \rightarrow \mathbb{R}$ is a function satisfying Assumption A.1 and $\mathbb{E}[F_n(Q)(i)] = f(Q)$ for any $i \in \mathcal{I}$ and $Q \in \mathbb{R}^{|\mathcal{I}|}$ where $f : \mathcal{I} \rightarrow \mathbb{R}$ is a function satisfying Assumption A.2.

Assumption A.1. 1) g is a max-norm non-expansion, 2) g is a span-norm non-expansion, 3) $g(x + ce) = g(x) + ce$ for any $c \in \mathbb{R}, x \in \mathbb{R}^{|\mathcal{I}|}$, 4) $g(cx) = cg(x)$ for any $c \in \mathbb{R}, x \in \mathbb{R}^{|\mathcal{I}|}$.

Assumption A.2. 1) f is L -Lipschitz, 2) there exists a positive scalar u s.t. $f(e) = u$ and $f(x + ce) = f(x) + cu$, 3) $f(cx) = cf(x)$.

Assumption A.3. For $n \in \{0, 1, 2, \dots\}$, $\mathbb{E}[\|R_n - r\|^2] \leq K$, $\mathbb{E}[\|G_n(Q) - g(Q)\|^2] \leq K(1 + \|Q\|^2)$ for any $Q \in \mathbb{R}^{|\mathcal{I}|}$, and $\mathbb{E}[\|F_n(Q) - f(Q)e\|^2] \leq K(1 + \|Q\|^2)$ for any $Q \in \mathbb{R}^{|\mathcal{I}|}$ for a suitable constant $K > 0$.

The above assumption means that the variances of $R_n, G_n(Q)$, and $F_n(Q)$ for any Q are bounded.

General RVI Q's update rule is

$$Q_{n+1}(i) \doteq Q_n(i) + \alpha_{\nu(n,i)}(R_n(i) - F_n(Q_n)(i) + G_n(Q_n)(i) - Q_n(i) + \epsilon_n(i))I\{i \in Y_n\}, \quad (\text{A.1})$$

where $\alpha_{\nu(n,i)}$ is the stepsize and ϵ_n is a sequence of random vectors of size $|\mathcal{I}|$.

We make following assumption on ϵ_n .

Assumption A.4 (Noise Assumption). $\|\epsilon_n\|_\infty \leq K(1 + \|Q_n\|_\infty)$ for some scalar K . Further, ϵ_n converges in probability to 0.

We make following assumptions on $\alpha_{\nu(n,i)}$.

Assumption A.5 (Stepsize Assumption). For all $n \geq 0$, $\alpha_n > 0$, $\sum_{n=0}^\infty \alpha_n = \infty$, and $\sum_{n=0}^\infty \alpha_n^2 < \infty$.

Assumption A.6 (Asynchronous Stepsize Assumption A). Let $[\cdot]$ denote the integer part of (\cdot) , for $x \in (0, 1)$,

$$\sup_i \frac{\alpha_{[xi]}}{\alpha_i} < \infty$$

and

$$\frac{\sum_{j=0}^{[yi]} \alpha_j}{\sum_{j=0}^i \alpha_j} \rightarrow 1$$

uniformly in $y \in [x, 1]$.

521 **Assumption A.7** (Asynchronous Stepsize Assumption B). *There exists $\Delta > 0$ such that*

$$\liminf_{n \rightarrow \infty} \frac{\nu(n, i)}{n+1} \geq \Delta,$$

522 *a.s., for all $s \in \mathcal{S}, o \in \mathcal{O}$. Furthermore, for all $x > 0$, let*

$$N(n, x) = \min \left\{ m > n : \sum_{i=n+1}^m \alpha_i \geq x \right\},$$

523 *the limit*

$$\lim_{n \rightarrow \infty} \frac{\sum_{i=\nu(n, i)}^{\nu(N(n, x), i)} \alpha_i}{\sum_{i=\nu(n, i')}^{\nu(N(n, x), i')} \alpha_i}$$

524 *exists a.s. for all s, s', o, o' .*

525 **Assumption A.8.** $r(i) - \bar{r} + g(q)(i) - q(i) = 0, \forall i \in \mathcal{I}$ has a unique solution for \bar{r} and a unique
526 *for q only up to a constant.*

527 Denoted the unique solution of \bar{r} by r_∞ . Further, it can be seen that the solution of q satisfying both
528 $r - \bar{r}e - g(q) - q = 0$ and $f(q) = r_\infty$ is unique because our assumption on f (Assumption A.2).
529 Denote the unique solution as q_∞ . We have,

$$f(q_\infty) = r_\infty. \quad (\text{A.2})$$

530 **Theorem A.1.** *Under Assumptions A.1-A.8, General RVI Q converges, almost surely, Q_n to q_∞ and*
531 *$f(Q_n)$ to r_∞ .*

532 *Proof.* Because (A.1) is in the same form as the asynchronous update (Equation 7.1.2) by Borkar
533 (2009), we apply the result in Section 7.4 of the same text (Borkar 2009) (see also Theorem 3.2
534 by Borkar (1998)) which shows convergence for Equation 7.1.2, to show the convergence of (A.1).
535 This result, given Assumption A.6 and A.7, only requires showing the convergence of the following
536 *synchronous* version of the General RVI Q algorithm:

$$Q_{n+1}(i) \doteq Q_n(i) + \alpha_n (R_n(i) - F_n(Q_n)(i) + g(Q_n)(i) - Q_n(i)) \quad \forall i \in \mathcal{I}. \quad (\text{A.3})$$

537 Define operators T_1, T_2 :

$$\begin{aligned} T_1(Q)(i) &\doteq r(i) + g(Q)(i) - r_\infty, \\ T_2(Q)(i) &\doteq r(i) + g(Q)(i) - f(Q) \\ &= T_1(Q)(i) + (r_\infty - f(Q)). \end{aligned}$$

538 Consider two ordinary differential equations (ODEs):

$$\dot{y}_t \doteq T_1(y_t) - y_t, \quad (\text{A.4})$$

$$\dot{x}_t \doteq T_2(x_t) - x_t = T_1(x_t) - x_t + (r_\infty - f(x_t))e. \quad (\text{A.5})$$

539 Note that because g is a non-expansion by Assumption A.1, both (A.4) and (A.5) have Lipschitz
540 R.H.S.'s and thus are well-posed.

541 Because g is a non-expansion, T_1 is also a non-expansion. Therefore we have the next lemma, which
542 restates Theorem 3.1 and Lemma 3.2 by Borkar and Soumyanath (1997).

543 **Lemma A.1.** *Let \bar{y} be an equilibrium point of (A.4). Then $\|y_t - \bar{y}\|_\infty$ is nonincreasing, and $y_t \rightarrow y_*$*
544 *for some equilibrium point y_* of (A.4) that may depend on y_0 .*

545 **Lemma A.2.** *(A.5) has a unique equilibrium at q_∞ .*

546 *Proof.* Because $f(q_\infty) = r_\infty$, we have that $q_\infty = T_1(q_\infty) = T_2(q_\infty)$, thus q_∞ is a equilibrium
547 point for (A.5). Conversely, if $T_2(Q) - Q = 0$, then $T_1Q + (r_\infty - f(Q))e = Q$. But the equation
548 $T_1Q + ce = Q$ only has a solution when $c = 0$ because of Assumption A.1. We have $c = 0$ and thus
549 $f(Q) = r_\infty$, which along with $T_1Q = Q$, implies $Q = q_\infty$. \square

550 **Lemma A.3.** Let $x_0 = y_0$, then $x_t = y_t + z_t e$, where z_t satisfies the ODE $\dot{z}_t = -uz_t + (r_\infty - f(y_t))$,
 551 and $k \doteq |\mathcal{I}|$.

552 *Proof.* From (A.4), (A.5), by the variation of parameters formula,

$$\begin{aligned} x_t &= \exp(-t)x_0 + \int_0^t \exp(\tau - t)T_1(x_\tau)d\tau + \left[\int_0^t \exp(\tau - t)(r_\infty - f(x_\tau))d\tau \right] e, \\ y_t &= \exp(-t)y_0 + \int_0^t \exp(\tau - t)T_1(y_\tau)d\tau. \end{aligned}$$

553 Then we have

$$\begin{aligned} &\max_{s,o}(x_t(s,o) - y_t(s,o)) \\ &\leq \int_0^t \exp(\tau - t) \max_{s,o}(T_1(x_\tau)(s,o) - T_1(y_\tau)(s,o))d\tau + \left[\int_0^t \exp(\tau - t)(r_\infty - f(x_\tau))d\tau \right] e, \\ &\min_{s,o}(x_t(s,o) - y_t(s,o)) \\ &\geq \int_0^t \exp(\tau - t) \min_{s,o}(T_1(x_\tau)(s,o) - T_1(y_\tau)(s,o))d\tau + \left[\int_0^t \exp(\tau - t)(r_\infty - f(x_\tau))d\tau \right] e. \end{aligned}$$

554 Subtracting, we have

$$sp(x_t - y_t) \leq \int_0^t \exp(\tau - t)sp(T_1(x_\tau) - T_1(y_\tau))d\tau,$$

555 where $sp(x)$ denotes the span of vector x .

556 Because we assumed that g is span-norm non-expansion, T_1 is also a span-norm non-expansion and
 557 thus

$$sp(x_t - y_t) \leq \int_0^t \exp(\tau - t)sp(T_1(x_\tau) - T_1(y_\tau))d\tau \leq \int_0^t \exp(\tau - t)sp(x_\tau - y_\tau)d\tau.$$

558 By Gronwall's inequality, $sp(x_t - y_t) = 0$ for all $t \geq 0$. Because $sp(x) = 0$ if and only if $x = ce$
 559 for some $c \in \mathbb{R}$, we have

$$x_t = y_t + z_t e, \quad t \geq 0.$$

560 for some z_t . Also $x_0 = y_0 \implies z_0 = 0$.

561 Now we show that $\dot{z}_t = -uz_t + (r_\infty - f(y_t))$. Note that $f(x_t) = f(y_t + z_t e) = f(y_t) + uz_t$. In
 562 addition, $T_1(x_t) - T_1(y_t) = T_1(y_t + z_t e) - T_1(y_t) = T_1(y_t) + z_t e - T_1(y_t) = z_t e$, therefore we
 563 have, for $z_t \in \mathbb{R}$:

$$\begin{aligned} \dot{z}_t e &= \dot{x}_t - \dot{y}_t \\ &= (T_1(x_t) - x_t + (r_\infty - f(x_t))e) - (T_1(y_t) - y_t) \quad (\text{from (A.4) and (A.5)}) \\ &= -(x_t - y_t) + (T_1(x_t) - T_1(y_t)) + (r_\infty - f(x_t))e \\ &= -z_t e + z_t e + (r_\infty - f(x_t))e \\ &= -uz_t e + uz_t e + (r_\infty - f(x_t))e \\ &= -uz_t e + (r_\infty - f(y_t))e \\ \implies \dot{z}_t &= -uz_t + (r_\infty - f(y_t)). \end{aligned}$$

564 □

565 **Lemma A.4.** q_∞ is the globally asymptotically stable equilibrium for (A.5).

566 *Proof.* We have shown that q_∞ is the unique equilibrium in Lemma A.2.

567 With that result, we first prove Lyapunov stability. That is, we need to show that given any $\epsilon > 0$, we
 568 can find a $\delta > 0$ such that $\|q_\infty - x_0\|_\infty \leq \delta$ implies $\|q_\infty - x_t\|_\infty \leq \epsilon$ for $t \geq 0$.

569 First, from Lemma A.3 we have $\dot{z}_t = -uz_t + (r_\infty - f(y_t))$. By variation of parameters and $z_0 = 0$,
 570 we have

$$z_t = \int_0^t \exp(u(\tau - t)) (r_\infty - f(y_\tau)) d\tau.$$

571 Then

$$\begin{aligned} \|q_\infty - x_t\|_\infty &= \|q_\infty - y_t - z_t u e\|_\infty \\ &\leq \|q_\infty - y_t\|_\infty + u |z_t| \\ &\leq \|q_\infty - y_0\|_\infty + u \int_0^t \exp(u(\tau - t)) |r_\infty - f(y_\tau)| d\tau \\ &= \|q_\infty - x_0\|_\infty + u \int_0^t \exp(u(\tau - t)) |f(q_\infty) - f(y_\tau)| d\tau \quad (\text{from (A.2)}). \end{aligned} \quad (\text{A.6})$$

572 Because f is L -lipschitz, we have

$$\begin{aligned} |f(q_\infty) - f(y_\tau)| &\leq L \|r_\infty - y_\tau\|_\infty \\ &\leq L \|r_\infty - y_0\|_\infty \quad (\text{from Lemma A.1}) \\ &= L \|r_\infty - x_0\|_\infty. \end{aligned}$$

573 Therefore

$$\begin{aligned} \int_0^t \exp(u(\tau - t)) |f(q_\infty) - f(y_\tau)| d\tau &\leq \int_0^t \exp(u(\tau - t)) L \|q_\infty - x_0\|_\infty d\tau \\ &= L \|q_\infty - x_0\|_\infty \int_0^t \exp(u(\tau - t)) d\tau \\ &= L \|q_\infty - x_0\|_\infty \frac{1}{u} (1 - \exp(-ut)) \\ &= \frac{L}{u} \|q_\infty - x_0\|_\infty (1 - \exp(-ut)). \end{aligned}$$

574 Substituting the above equation in (A.6), we have

$$\|q_\infty - x_t\|_\infty \leq (1 + L) \|q_\infty - x_0\|_\infty.$$

575 Lyapunov stability follows.

576 Now in order to prove the asymptotic stability, in addition to Lyapunov stability, we need to show
 577 that there exists $\delta > 0$ such that if $\|x_0 - q_\infty\|_\infty < \delta$, then $\lim_{t \rightarrow \infty} \|x_t - q_\infty\|_\infty = 0$. Note that

$$\begin{aligned} \lim_{t \rightarrow \infty} z_t &= \lim_{t \rightarrow \infty} \int_0^t \exp(u(\tau - t)) (r_\infty - f(y_\tau)) d\tau \\ &= \lim_{t \rightarrow \infty} \frac{\int_0^t \exp(u\tau) (r_\infty - f(y_\tau)) d\tau}{\exp(ut)} \\ &= \lim_{t \rightarrow \infty} \frac{\exp(ut) (r_\infty - f(y_t))}{u \exp(ut)} \quad (\text{by L'Hospital's rule}) \\ &= \frac{r_\infty - f(y_\infty)}{u} \quad (\text{by Lemma A.1}). \end{aligned}$$

578 Because $x_t = y_t + z_t e$ (Lemma A.3) and $y_t \rightarrow y_\infty$ (Lemma A.1), we have $x_t \rightarrow y_\infty + (r_\infty -$
 579 $f(y_\infty))e/u$, which must coincide with q_∞ because that is the only equilibrium point for (A.5)
 580 (Lemma A.2). Therefore $\lim_{t \rightarrow \infty} \|x_t - q_\infty\|_\infty = 0$ for any x_0 . Asymptotic stability is shown and
 581 the proof is complete. \square

582 **Lemma A.5.** Equation A.3 converges a.s. Q_n to q_∞ as $n \rightarrow \infty$.

583 *Proof.* The proof uses Borkar's (2008) Theorem 2 in Section 2 and is essentially the same as Lemma
 584 3.8 by Abounadi et al. (2001). For completeness, we repeat the proof (with more details) here.

585 First write the synchronous update (A.3) as

$$Q_{n+1} = Q_n + \alpha_n(h(Q_n) + M_{n+1} + \epsilon_n),$$

586 where

$$\begin{aligned} h(Q_n)(i) &\doteq r(i) - f(Q_n) + g(Q_n)(i) - Q_n(i) \\ &= T_2(Q_n)(i) - Q_n(i), \\ M_{n+1}(i) &\doteq R_n(i) - F_n(Q_n)(i) + G_n(Q_n)(i) - T_2(Q_n)(i). \end{aligned}$$

587 It can be shown that ϵ_n is asymptotically negligible and therefore does not affect the conclusions of
 588 Theorem 2 (text after Equation B.66 by Wan et al. 2021).

589 Theorem 2 requires verifying following conditions and concludes that Q_n converges to a (possibly
 590 sample path dependent) compact connected internally chain transitive invariant set of ODE $\dot{x}_t =$
 591 $h(x_t)$. This is exactly the ODE defined in (A.5). Lemma A.2 and A.4 conclude that this ODE has
 592 q_∞ as the unique globally asymptotically stable equilibrium. Therefore the (possibly sample path
 593 dependent) compact connected internally chain transitive invariant set is a singleton set containing
 594 only the unique globally asymptotically stable equilibrium. Thus Theorem 2 concludes that $Q_n \rightarrow q_\infty$
 595 a.s. as $n \rightarrow \infty$. We now list conditions required by Theorem 2:

- 596 • **(A1)** The function h is Lipschitz: $\|h(x) - h(y)\| \leq L \|x - y\|$ for some $0 < L < \infty$.
- 597 • **(A2)** The sequence $\{\alpha_n\}$ satisfies $\alpha_n > 0$, and $\sum \alpha_n = \infty$, $\sum \alpha_n^2 < \infty$.
- 598 • **(A3)** $\{M_n\}$ is a martingale difference sequence with respect to the increasing family of
 599 σ -fields

$$\mathcal{F}_n \doteq \sigma(Q_i, M_i, i \leq n), n \geq 0.$$

600 That is

$$\mathbb{E}[M_{n+1} \mid \mathcal{F}_n] = 0 \quad \text{a.s., } n \geq 0.$$

601 Furthermore, $\{M_n\}$ are square-integrable

$$\mathbb{E}[\|M_{n+1}\|^2 \mid \mathcal{F}_n] \leq K(1 + \|Q_n\|^2) \quad \text{a.s., } n \geq 0,$$

602 for some constant $K > 0$.

- 603 • **(A4)** $\sup_n \|Q_n\| \leq \infty$ a.s..

604 Let us verify these conditions now.

605 (A1) is satisfied because T_2 is Lipschitz.

606 (A2) is satisfied by Assumption A.5.

607 (A3) is also satisfied because for any $i \in \mathcal{I}$

$$\begin{aligned} \mathbb{E}[M_{n+1}(i) \mid \mathcal{F}_n] &= \mathbb{E}[R_n(i) - F_n(Q_n)(i) + G_n(i) - T_2(Q_n)(i) \mid \mathcal{F}_n] \\ &= \mathbb{E}[R_n(i) - F_n(Q_n)(i) + G_n(Q_n)(i) \mid \mathcal{F}_n] - T_2(Q_n)(i) \\ &= 0, \end{aligned}$$

608 and $\mathbb{E}[\|M_{n+1}\|^2 \mid \mathcal{F}_n] \leq \mathbb{E}[\|R_n - r\|^2 \mid \mathcal{F}_n] + \mathbb{E}[\|F_n(Q_n) - f(Q_n)e\|^2 \mid \mathcal{F}_n] +$
 609 $\mathbb{E}[\|G_n(Q_n) - g(Q_n)\|^2 \mid \mathcal{F}_n] \leq K(1 + \|Q_n\|^2)$ for a suitable constant $K > 0$ can be verified
 610 by a simple application of triangle inequality.

611 To verify (A4), we apply Theorem 7 in Section 3 by Borkar (2008), which shows $\sup_n \|Q_n\| \leq \infty$
 612 a.s., if (A1), (A2), and (A3) are all satisfied and in addition we have the following condition satisfied:

613 **(A5)** The functions $h_d(x) \doteq h(dx)/d$, $d \geq 1$, $x \in \mathbb{R}^k$, satisfy $h_d(x) \rightarrow h_\infty(x)$ as $d \rightarrow \infty$, uniformly
 614 on compacts for some $h_\infty \in C(\mathbb{R}^k)$. Furthermore, the ODE $\dot{x}_t = h_\infty(x_t)$ has the origin as its unique
 615 globally asymptotically stable equilibrium.

616 Note that

$$h_\infty(x) = \lim_{d \rightarrow \infty} h_d(x) = \lim_{d \rightarrow \infty} (T_2(dx) - dx) / d = g(x) - f(x)e - x,$$

617 because $g(cx) = cg(x)$ and $f(cx) = cf(x)$ by our assumption.

618 The function h_∞ is clearly continuous in every $x \in \mathbb{R}^k$ and therefore $h_\infty \in C(\mathbb{R}^k)$.

619 Now consider the ODE $\dot{x}_t = h_\infty(x_t) = g(x_t) - f(x_t)e - x_t$. Clearly the origin is an equilibrium.
 620 This ODE is a special case of (A.5), corresponding to the $r(s, o) \forall s \in \mathcal{S}, o \in \mathcal{O}$ being always
 621 zero. Therefore Lemma A.2 and A.4 also apply to this ODE and the origin is the unique globally
 622 asymptotically stable equilibrium.

623 (A1), (A2), (A3), (A4) are all verified and therefore

$$Q_n \rightarrow q_\infty \text{ a.s. as } n \rightarrow \infty.$$

624

□

625

□

626 A.2 Theorem 1

627 For simplicity, we will only provide formal theorems and proofs for our *control* learning and planning
 628 algorithms. The formal theorems and proofs for our prediction algorithms are similar to those for
 629 the control algorithms and are thus omitted. To this end, we first provide a general algorithm that
 630 includes both learning and planning control algorithms. We call it *General Inter-option Differential*
 631 *Q*. We first formally define it and then explain why both inter-option Differential Q-learning and
 632 inter-option Differential Q-planning are special cases of General Inter-option Differential Q. We then
 633 provide assumptions and the convergence theorem of the general algorithm. The theorem would lead
 634 to convergence of the special cases. Finally, we provide a proof for the theorem.

635 Given an SMDP $\hat{\mathcal{M}} = (\mathcal{S}, \mathcal{O}, \hat{\mathcal{L}}, \hat{\mathcal{R}}, \hat{p})$, for each state $s \in \mathcal{S}$, option $o \in \mathcal{O}$, and discrete step $n \geq 0$,
 636 let $\hat{R}_n(s, o), \hat{S}'_n(s, o), \hat{L}_n(s, o) \sim \hat{p}(\cdot, \cdot, \cdot | s, o)$ denote a sample of resulting state, reward and the
 637 length. We hypothesize a set-valued process $\{Y_n\}$ taking values in the set of nonempty subsets of
 638 $\mathcal{S} \times \mathcal{O}$ with the interpretation: $Y_n = \{(s, o) : (s, o) \text{ component of } Q \text{ was updated at time } n\}$. Let
 639 $\nu(n, s, o) \doteq \sum_{k=0}^n I\{(s, o) \in Y_k\}$, where I is the indicator function. Thus $\nu(n, s, o)$ = the number
 640 of times the (s, o) component was updated up to step n . The update rules of General Inter-option
 641 Differential Q are

$$Q_{n+1}(s, o) \doteq Q_n(s, o) + \alpha_{\nu(n, s, o)} \delta_n(s, o) / L_n(s, o) I\{(s, o) \in Y_n\}, \quad \forall s \in \mathcal{S}, o \in \mathcal{O}, \quad (\text{A.7})$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta \sum_{s, o} \alpha_{\nu(n, s, o)} \delta_n(s, o) / L_n(s, o) I\{(s, o) \in Y_n\}, \quad (\text{A.8})$$

$$L_{n+1}(s, o) \doteq L_n(s, o) + \beta_n(s, o) (\hat{L}_n(s, o) - L_n(s, o)) I\{(s, o) \in Y_n\}, \quad (\text{A.9})$$

642 where

$$\delta_n(s, o) \doteq \hat{R}_n(s, o) - \bar{R}_n L_n(s, o) + \max_{o'} Q_n(\hat{S}'_n(s, o), o') - Q_n(s, o) \quad (\text{A.10})$$

643 is the TD error.

644 Here $\alpha_{\nu(n, s, o)}$ is the stepsize at step n for state-action pair (s, o) . The quantity $\alpha_{\nu(n, s, o)}$ depends
 645 on the sequence $\{\alpha_n\}$, which is an algorithmic design choice, and also depends on the visitation
 646 of state-option pairs $\nu(n, s, o)$. To obtain the stepsize, the algorithm could maintain a $|\mathcal{S} \times \mathcal{O}|$ -size
 647 table counting the number of visitations to each state-option pair, which is exactly $\nu(\cdot, \cdot, \cdot)$. Then the
 648 stepsize $\alpha_{\nu(n, s, o)}$ can be obtained as long as the sequence $\{\alpha_n\}$ is specified.

649 Q_0 and R_0 can be initialized arbitrarily. Note that L_0 can not be initialized to 0 because it is the
 650 divisor for both (A.7) and (A.8) for the first update. Because the expected length of all options would
 651 be greater than or equal to 1, we choose L_0 to be 1. In this way, L_n will never be 0 because it
 652 is initialized to 1 and all the sampled option lengths are greater than or equal to 1. Therefore the
 653 problem of division by 0 will not happen in the updates.

654 Now we show inter-option Differential Q-learning and inter-option Differential Q-planning are
 655 special cases of General Inter-option Differential Q. Consider a sequence of real experience
 656 $\dots, \hat{S}_n, \hat{O}_n, \hat{R}_n, \hat{L}_n, \hat{S}_{n+1}, \dots$

$$\begin{aligned} Y_n(s, o) &= 1, \text{ if } s = \hat{S}_n, o = \hat{O}_n, \\ Y_n(s, o) &= 0 \text{ otherwise,} \end{aligned}$$

657 and $\hat{S}'_n(\hat{S}_n, \hat{O}_n) = \hat{S}_{n+1}$, $\hat{R}_n(\hat{S}_n, \hat{O}_n) = \hat{R}_{n+1}$, $\hat{L}_n(\hat{S}_n, \hat{O}_n) = \hat{L}_n$, update rules (A.7), (A.8), and
 658 (A.10) become

$$\begin{aligned} Q_{n+1}(\hat{S}_n, \hat{O}_n) &\doteq Q_n(\hat{S}_n, \hat{O}_n) + \alpha_{\nu(n, \hat{S}_n, \hat{O}_n)} \hat{\delta}_n / L_n(\hat{S}_n, \hat{O}_n), \text{ and } Q_{n+1}(s, o) \doteq Q_n(s, o), \forall s \neq \hat{S}_n, o \neq \hat{O}_n, \\ \bar{R}_{n+1} &\doteq \bar{R}_n + \eta \alpha_{\nu(n, \hat{S}_n, \hat{O}_n)} \hat{\delta}_n / L_n(\hat{S}_n, \hat{O}_n), \\ \hat{\delta}_n &\doteq \hat{R}_n - \bar{R}_n \hat{L}_n + \max_{o'} Q_n(\hat{S}_{n+1}, o') - Q_n(\hat{S}_n, \hat{O}_n), \\ L_{n+1}(\hat{S}_n, \hat{O}_n) &\doteq L_n(\hat{S}_n, \hat{O}_n) + \beta_n(\hat{S}_n, \hat{O}_n)(\hat{L}_n - L_n(\hat{S}_n, \hat{O}_n)) \end{aligned}$$

659 which are inter-option Differential Q-learning's update rules (Section 3) with stepsize α in the n -th
 660 update being $\alpha_{\nu(n, \hat{S}_n, \hat{O}_n)}$, and the stepsize β being $\beta(\hat{S}_n, \hat{O}_n)$.

661 If we consider a sequence of simulated experience $\dots, \tilde{S}_n, \tilde{O}_n, \tilde{R}_n, \tilde{L}_n, \tilde{S}'_n, \dots$

$$\begin{aligned} Y_n(s, o) &= 1, \text{ if } s = \tilde{S}_n, o = \tilde{O}_n, \\ Y_n(s, o) &= 0 \text{ otherwise,} \end{aligned}$$

662 and $\tilde{S}'_n(s, o) = \tilde{S}'_n$, $\tilde{R}_n(s, o) = \tilde{R}_n$, $\tilde{L}_n(s, o) = \tilde{L}_n$, update rules (A.7)-(A.10) become

$$\begin{aligned} Q_{n+1}(\tilde{S}_n, \tilde{O}_n) &\doteq Q_n(\tilde{S}_n, \tilde{O}_n) + \alpha_{\nu(n, \tilde{S}_n, \tilde{O}_n)} \tilde{\delta}_n / L_n, \text{ and } Q_{n+1}(s, o) \doteq Q_n(s, o), \forall s \neq \tilde{S}_n, o \neq \tilde{O}_n, \\ \bar{R}_{n+1} &\doteq \bar{R}_n + \eta \alpha_{\nu(n, \tilde{S}_n, \tilde{O}_n)} \tilde{\delta}_n / L_n, \\ \tilde{\delta}_n &\doteq \tilde{R}_n - \bar{R}_n \tilde{L}_n + \max_{o'} Q_n(\tilde{S}'_n, o') - Q_n(\tilde{S}_n, \tilde{O}_n), \\ L_{n+1}(\tilde{S}_n, \tilde{O}_n) &\doteq L_n(\tilde{S}_n, \tilde{O}_n) + \beta_n(\tilde{S}_n, \tilde{O}_n)(\tilde{L}_n - L_n(\tilde{S}_n, \tilde{O}_n)). \end{aligned}$$

663 Now, in the planning setting, the model can produce an expected length, instead of a sampled one.
 664 And there estimating the expected length using L_n is no longer needed. The above updates reduce to

$$\begin{aligned} Q_{n+1}(\tilde{S}_n, \tilde{O}_n) &\doteq Q_n(\tilde{S}_n, \tilde{O}_n) + \alpha_{\nu(n, \tilde{S}_n, \tilde{O}_n)} \tilde{\delta}_n / \tilde{L}_n, \text{ and } Q_{n+1}(s, o) \doteq Q_n(s, o), \forall s \neq \tilde{S}_n, o \neq \tilde{O}_n, \\ \bar{R}_{n+1} &\doteq \bar{R}_n + \eta \alpha_{\nu(n, \tilde{S}_n, \tilde{O}_n)} \tilde{\delta}_n / \tilde{L}_n, \\ \tilde{\delta}_n &\doteq \tilde{R}_n - \bar{R}_n \tilde{L}_n + \max_{o'} Q_n(\tilde{S}'_n, o') - Q_n(\tilde{S}_n, \tilde{O}_n). \end{aligned}$$

665 The above update rules are our inter-option Differential Q-planning's update rules with stepsize at
 666 planning step n being $\alpha_{\nu(n, \tilde{S}_n, \tilde{O}_n)}$.

667 We now provide a theorem, along with its proof, showing the convergence of General Inter-option
 668 Differential Q.

669 **Theorem A.2.** Under Assumptions 1, A.5, A.6, A.7, and that $0 \leq \beta_n(s, o) \leq 1$, $\sum_n \beta_n(s, o) = \infty$,
 670 and $\sum_n \beta_n^2(s, o) < \infty$, and $\beta_n(s, o) = 0$ unless $s = \hat{S}_n$, General Inter-option Differential Q
 671 (Equations A.7-A.10) converges, almost surely, Q_n to q satisfying both (2) and

$$\eta(\sum q - \sum Q_0) = r_* - \bar{R}_0,$$

672 \bar{R}_n to r_* , and $r(\mu_n)$ to r_* where μ_n is a greedy policy w.r.t. Q_n .

673 *Proof.* At each step, the increment to \bar{R}_n is η times the increment to Q_n and $\sum Q_n$. Therefore, the
 674 cumulative increment can be written

$$\begin{aligned}\bar{R}_n - \bar{R}_0 &= \eta \sum_{i=0}^{n-1} \sum_{s,o} \alpha_{\nu(i,s,o)} \delta_i(s,o) / L_i(s,o) I\{(s,o) \in Y_i\} \\ &= \eta \left(\sum Q_n - \sum Q_0 \right) \\ \implies \bar{R}_n &= \eta \sum Q_n - \eta \sum Q_0 + \bar{R}_0 = \eta \sum Q_n - c, \tag{A.11}\end{aligned}$$

$$\text{where } c \doteq \eta \sum Q_0 - \bar{R}_0. \tag{A.12}$$

675 Now substituting \bar{R}_n in (A.7) with (A.11), we have $\forall s \in \mathcal{S}, o \in \mathcal{O}$:

$$\begin{aligned}Q_{n+1}(s,o) &= Q_n(s,o) + \alpha_{\nu(n,s,o)} \\ \frac{\hat{R}_n(s,o) - L_n(s,o)(\eta \sum Q_n - c) + \max_{o'} Q_n(\hat{S}'_n(s,o), o') - Q_n(s,o)}{L_n(s,o)} I\{(s,o) \in Y_n\} \\ &= Q_n(s,o) + \alpha_{\nu(n,s,o)} \\ &\quad \left(\frac{\hat{R}_n(s,o) - l_n(s,o)(\eta \sum Q_n - c) + \max_{o'} Q_n(\hat{S}'_n(s,o), o') - Q_n(s,o)}{l(s,o)} + \epsilon_n(s,o) \right) I\{(s,o) \in Y_n\}, \tag{A.13}\end{aligned}$$

676 where $l(s,o)$ is the expected length of option o , starting from state s , and $\epsilon_n(s,o) \doteq (\hat{R}_n(s,o) -$
 677 $L_n(s,o)(\eta \sum Q_n - c) + \max_{o'} Q_n(\hat{S}'_n(s,o), o') - Q_n(s,o)) / L(s,o) - (\hat{R}_n(s,o) - l(s,o)(\eta \sum Q_n -$
 678 $c) + \max_{o'} Q_n(\hat{S}'_n(s,o), o') - Q_n(s,o)) / l(s,o)$.

679 Standard stochastic approximation result can be applied to show that L_n converges to l . Further, it
 680 can be seen that $\|\epsilon_n\|_\infty \leq K(1 + \|Q_n\|)$ for some positive K and, by continuous
 681 mapping theorem, converges to 0 almost surely (and thus in probability).

682 We now show that (A.13) is a special case of (A.1). To see this point, let

$$\begin{aligned}i &= (s,o), \\ R_n(i) &= \frac{\hat{R}_n(s,o)}{l(s,o)} + c, \\ G_n(Q_n)(i) &= \frac{\max_{o'} Q_n(\hat{S}'_n(s,o), o')}{l(s,o)} + \frac{l(s,o) - 1}{l(s,o)} Q_n(s,o), \\ F(Q_n)(i) &= \eta \sum Q_n, \\ \epsilon_n(i) &= \epsilon_n(s,o).\end{aligned}$$

683 We now verify the assumptions of Theorem A.1 for Inter-option General Differential Q. Assump-
 684 tion A.1 and Assumption A.2 can be verified easily. Assumption A.3 satisfies because the MDP
 685 is finite. Assumption A.4 is satisfied as shown above. Assumption A.5-A.7 are satisfied due to
 686 assumptions of the theorem being proved. Assumption A.8 is satisfied because

$$\begin{aligned}r(i) - \bar{r} + g(q)(i) - q(i) &= \mathbb{E}[R_n(i) - \bar{r} + G_n(q)(i) - q(i)] \\ &= \mathbb{E} \left[\frac{\hat{R}_n(s,o) + cl(s,o) - \bar{r}l(s,o) + \max_{o'} q(\hat{S}'_n(s,o), o') + (l(s,o) - 1)q(s,o) - l(s,o)q(s,o)}{l(s,o)} \right] \\ &= \frac{\mathbb{E} \left[\hat{R}_n(s,o) + cl(s,o) - \bar{r}l(s,o) + \max_{o'} q(\hat{S}'_n(s,o), o') - q(s,o) \right]}{l(s,o)}.\end{aligned}$$

687 From (2) we know if the above equation equals to 0, then under Assumption 1, $\bar{r} = r_* + c$ is the
 688 unique solution and the solutions for q form a set $q = q_* + ce$.

689 All the assumptions are verified and thus from Theorem A.1 we conclude that Q_n converges to a point
 690 satisfying $\eta \sum q = r_* + c = r_* + \eta \sum Q_0 - \bar{R}_0$ and $\bar{R}_n = \eta \sum Q_n - c$ to $\eta \sum q - c = r_* + c - c = r_*$.
 691 Finally, in order to show $r(\mu_n) \rightarrow r_*$, we first extend Theorem 8.5.5 by Puterman (1994) to deal with
 692 SMDP.

693 **Lemma A.6.** *Under Assumption 1, $\forall Q \in \mathbb{R}^{|\mathcal{S} \times \mathcal{O}|}$*

$$\min_{s,o} TQ(s, o) \leq r(\mu_Q) \leq r_* \leq \max_{s,o} TQ(s, o),$$

694 where $TQ(s, o) \doteq \sum_{s',r,l} \hat{p}(s', r, l | s, o)(r + \max_{o'} Q(s', o'))$ and μ_Q denotes a greedy policy w.r.t.
 695 Q .

696 *Proof.* Note that

$$r(\mu_Q) = \sum_{s',r,l} \hat{p}(s', r, l | s, o)(r + \sum_{o'} \mu_Q(o' | s')Q(s', o') - Q(s, o)).$$

697 Therefore

$$\begin{aligned} \min_{s,o} (TQ_n(s, o) - Q_n(s, o)) &\leq r(\mu_n) \leq r_* \leq \max_{s,o} (TQ_n(s, o) - Q_n(s, o)) \\ &\implies |r_* - r(\mu_n)| \leq sp(TQ_n - Q_n). \end{aligned}$$

698 □

699 Because $Q_n \rightarrow q_\infty$ a.s., and $sp(TQ_n - Q_n)$ is a continuous function of Q_n , by continuous mapping
 700 theorem, $sp(TQ_n - Q_n) \rightarrow sp(Tq_\infty - q_\infty) = 0$ a.s. Therefore we conclude that $r(\mu_n) \rightarrow r_*$.

701 □

702 A.3 Theorem 2

703 The proof for the intra-option evaluation equation is simple. First note that these equations can be
 704 written in the vector form:

$$0 = r - \bar{r}e + (P_\mu - I)q,$$

705 where $r(s, o) = \mathbb{E}[R_{t+1} | S_t = s, O_t = o]$, $P_\mu((s, o), (s', o')) \doteq \Pr(S_{t+1} = s', O_{t+1} = o' | S_t =$
 706 $s, O_t = o, \mu) = \beta(s', o)\mu(o' | s') + (1 - \beta(s', o))\mathbb{I}(o = o')$, and e is a all-one vector. Intuitively, the
 707 intra-option evaluation equation can be viewed as the evaluation equation for some average-reward
 708 MRP with reward and dynamics being defined as r and P_μ .

709 By Theorem 8.2.6 and Corollary 8.2.7 in Puterman (1994), the intra-option evaluation equation part
 710 in Theorem 2 is shown as long as the Markov chain associated with P_μ is unichain. Note that by
 711 our Assumption 1, there is only one recurrent class of states under any policy. Therefore no matter
 712 what the start state-option pair is, the agent will enter in the same recurrent class of states. Therefore
 713 we have, for every state \bar{s} in the recurrent class and an option \bar{o} such that $\mu(\bar{o} | \bar{s}) > 0$, the MDP
 714 visits (\bar{s}, \bar{o}) an infinite number of times. This shows that any two state-option pairs can not be in
 715 two separate recurrent sets of state-option pairs. Therefore the Markov chain associated with P_μ is
 716 unichain.

717 The proof for the Intra-option Optimality Equations is more complicated. First, similar as what we
 718 know in the discounted primitive action case, we have the following lemma for the discounted option
 719 case.

720 **Lemma A.7.** *For every $0 < \gamma < 1$, there exists a stationary deterministic discount optimal policy.*

721 The proof uses similar arguments as Theorem 6.2.10 and Proposition 4.4.3 by Puterman (1994).

722 Now choose a sequence of discount factors $\{\gamma_n\}$, $0 \leq \gamma_n < 1$ with the property that $\gamma_n \uparrow 1$. By
 723 lemma A.7, for each γ_n , there exists a stationary discount optimal policy. Because the total number
 724 of Markov deterministic policies is finite, we can choose a subsequence $\{\gamma'_n\}$ for which the same
 725 Markov deterministic policy, μ , is discount optimal for all γ'_n . Denote this subsequence by $\{\gamma_n\}$.

726 Because μ is discount optimal for $\gamma_n, \forall n, q_*^{\gamma_n} = q_\mu^{\gamma_n}, \forall n$. By intra-option optimality equations in the
 727 discounted case (Sutton et al., 1999), for all $s \in \mathcal{S}, o \in \mathcal{O}$,

$$\begin{aligned} 0 &= \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) \left(r + \gamma_n \beta(s', o) q_\mu^{\gamma_n}(s', \mu(s')) + \gamma_n (1 - \beta(s', o)) q_\mu^{\gamma_n}(s', o) \right) - q_\mu^{\gamma_n}(s, o) \\ &= \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) \left(r + \gamma_n \beta(s', o) \max_{o'} q_\mu^{\gamma_n}(s', o') + \gamma_n (1 - \beta(s', o)) q_\mu^{\gamma_n}(s', o) \right) - q_\mu^{\gamma_n}(s, o). \end{aligned} \quad (\text{A.14})$$

728 By corollary 8.2.4 by Puterman (1994),

$$q_\mu^{\gamma_n} = (1 - \gamma_n)^{-1} r(\mu) + q_\mu + f(\gamma_n), \quad (\text{A.15})$$

729 where $r(\mu)$ and q_μ are the reward rate and differential value function under policy μ , and $f(\gamma)$ is a
 730 function of γ that converges to 0 as $\gamma \uparrow 1$.

731 Substituting (A.15) into (A.14), we have

$$\begin{aligned} 0 &= \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) (r + \gamma_n \beta(s', o) \max_{o'} [(1 - \gamma_n)^{-1} r(\mu) + q_\mu(s', o') + f(\gamma_n, s', o')]) \\ &\quad + \gamma_n (1 - \beta(s', o)) [(1 - \gamma_n)^{-1} r(\mu) + q_\mu(s', o) + f(\gamma_n, s', o)] \\ &\quad - [(1 - \gamma_n)^{-1} r(\mu) + q_\mu(s, o) + f(\gamma_n, s, o)] \\ &= \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) (r - r(\mu) + \gamma_n \beta(s', o) \max_{o'} [q_\mu(s', o') + f(\gamma_n, s', o')]) \\ &\quad + \gamma_n (1 - \beta(s', o)) [q_\mu(s', o) + f(\gamma_n, s', o)] \\ &\quad - [q_\mu(s, o) + f(\gamma_n, s, o)] \\ &= \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) (r - r(\mu) + \beta(s', o) \max_{o'} [q_\mu(s', o') + f(\gamma_n, s', o')]) \\ &\quad + (\gamma_n - 1) \beta(s', o) \max_{o'} [q_\mu(s', o') + f(\gamma_n, s', o')] \\ &\quad + (1 - \beta(s', o)) [q_\mu(s', o) + f(\gamma_n, s', o)] \\ &\quad + (\gamma_n - 1) (1 - \beta(s', o)) [q_\mu(s', o) + f(\gamma_n, s', o)] \\ &\quad - [q_\mu(s, o) + f(\gamma_n, s, o)]. \end{aligned}$$

732 Note that $(\gamma - 1) \beta(s', o) \max_{o'} [q_\mu(s', o') + f(\gamma, s', o')]$ and $(\gamma - 1) (1 - \beta(s', o)) [q_\mu(s', o) +$
 733 $f(\gamma, s', o)]$ both converge to 0 as $\gamma \uparrow 1$.

734 Now take $n \rightarrow \infty$, then $\gamma_n \uparrow 1$, we have

$$0 = \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) \left(r - r(\mu) + \beta(s', o) \max_{o'} q_\mu(s', o') + (1 - \beta(s', o)) q_\mu(s', o) \right) - q_\mu(s, o).$$

735 We see that $\bar{r} = r(\mu)$ and $q = q_\mu$ is a solution of (10)-(11).

736 Now we show that the solution for \bar{r} is unique. Define

$$B(\bar{r}, q) \doteq \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) \left(r - \bar{r} + \beta(s', o) \max_{o'} q(s', o') + (1 - \beta(s', o)) q(s', o) \right) - q(s, o).$$

737 First we show if $B(\bar{r}, q) = 0$, then $\bar{r} \geq r_*$.

$$\begin{aligned} 0 &= B(\bar{r}, q) \\ &= \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) (r - \bar{r} + \beta(s', o) \max_{o'} q(s', o') + (1 - \beta(s', o)) q(s', o)) - q(s, o) \\ &\geq \sup_{\mu \in \Pi^{MR}} \sum_a \pi(a|s, o) \sum_{s', r} p(s', r|s, a) \\ &\quad \left(r - \bar{r} + \beta(s', o) \sum_{o'} \mu(o'|s') q(s', o') + (1 - \beta(s', o)) q(s', o) \right) - q(s, o), \end{aligned}$$

738 where Π^{MR} denotes the set of all Markov randomized policies. In vector form, the above equation
 739 can be written as:

$$0 \geq \sup_{\mu \in \Pi^{MR}} \{r - \bar{r}e + (P_\mu - I)q\}.$$

740 Therefore $\forall \mu \in \Pi^{MR}$,

$$\bar{r}e \geq r + (P_\mu - I)q.$$

741 Apply P_μ to both sides,

$$\begin{aligned} P_\mu \bar{r}e &\geq P_\mu r + P_\mu (P_\mu - I)q, \\ \bar{r}e &\geq P_\mu r + P_\mu (P_\mu - I)q. \end{aligned}$$

742 Repeating this process we have:

$$\bar{r}e \geq P_\mu^n r + P_\mu^n (P_\mu - I)q.$$

743 Summing these expressions from $n = 0$ to $n = N - 1$ we have:

$$N\bar{r}e \geq \sum_{n=0}^{N-1} (P_\mu^n r + P_\mu^n (P_\mu - I)q) = \sum_{n=0}^{N-1} P_\mu^n r + (P_\mu^N - P_\mu^{N-1})q.$$

744 Because $\lim_{N \rightarrow \infty} \frac{1}{N} (P_\mu^N - P_\mu^{N-1})q = 0$,

$$\bar{r}e \geq \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} P_\mu^n r = r(\mu)e,$$

745 for all $\mu \in \Pi^{MR}$. Therefore $\bar{r} \geq r_*$.

746 Then we show that if $0 = B(\bar{r}, q)$ then $\bar{r} \leq r_*$. As we proved above, if (\bar{r}, q) satisfies that $0 = B(\bar{r}, q)$
 747 then there exists a policy μ such that $\bar{r}e = r + (P_\mu - I)q$ is true. Therefore,

$$\begin{aligned} P_\mu^n \bar{r}e &= P_\mu^n r + P_\mu^n (P_\mu - I)q, \\ \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} P_\mu^n \bar{r}e &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} (P_\mu^n r + P_\mu^n (P_\mu - I)q), \\ \bar{r}e &= \lim_{N \rightarrow \infty} \sum_{n=0}^{N-1} P_\mu^n r = r(\mu)e \leq r_*e. \end{aligned}$$

748 Therefore $\bar{r} \leq r_*$. Combining $\bar{r} \geq r_*$ and $\bar{r} \leq r_*$ we have $\bar{r} = r_*$.

749 Finally, we show that the solution for q is unique only up to a constant. Note that one could iteratively
 750 replace q in the r.h.s. of the intra-option Optimality equation (10)-(11) by the entire r.h.s. of the
 751 intra-option Optimality equation, resulting to the inter-option Optimality equation (2). Therefore any
 752 solution of (10)-(11) must be a solution of (2). But we know that the solutions for q in (2) is unique
 753 only up to a constant. Therefore the solutions for q in (10)-(11) can not differ by a non-constant.
 754 Further, it is easy to see that if q is a solution, then $q + ce, \forall c$ is also a solution. The theorem is
 755 proved.

756 \square

757 A.4 Theorem 3

758 For simplicity, we will only provide formal theorems and proofs for our *control* learning and planning
 759 algorithms. The formal theorems and proofs for our prediction algorithms are similar to those for
 760 the control algorithms and are thus omitted. To this end, we first provide a general algorithm that
 761 includes both learning and planning control algorithms. We call it *General Intra-option Differential*
 762 *Q*. We first formally define it and then explain why both Intra-option Differential Q-learning and
 763 Intra-option Differential Q-planning are special cases of General Intra-option Differential Algorithm.
 764 We then provide assumptions and the convergence theorem of the general algorithm. The theorem
 765 would lead to convergence of the special cases. Finally, we provide a proof for the theorem.

Given an MDP $\mathcal{M} \doteq (\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$ and a set of options \mathcal{O} , for each state $s \in \mathcal{S}$, option $o \in \mathcal{O}$, a reference option \bar{o} , and discrete step $n \geq 0$, let $A_n(s, \bar{o}) \sim \pi(\cdot \mid s, \bar{o})$, $R_n(s, A_n(s, \bar{o}))$, $S'_n(s, A_n(s, \bar{o})) \sim p(\cdot, \cdot \mid s, A_n(s, \bar{o}))$ denote, given state-option pair (s, \bar{o}) , a sample of the chosen action and the resulting state and reward. We hypothesize a set-valued process $\{Y_n\}$ taking values in the set of nonempty subsets of $\mathcal{S} \times \mathcal{O}$ with the interpretation: $Y_n = \{(s, o) : (s, o) \text{ component of } Q \text{ was updated at time } n\}$. Let $\nu(n, s, o) \doteq \sum_{k=0}^n I\{(s, o) \in Y_k\}$, where I is the indicator function. Thus $\nu(n, s, o) =$ the number of times the (s, o) component was updated up to step n . In addition, we hypothesize a set-valued process $\{Z_n\}$ taking values in the set of nonempty subsets of \mathcal{O} with the interpretation: $Z_n = \{\bar{o} : \bar{o} \text{ component was visited at time } n\}$. $\sum_{\bar{o}} I\{\bar{o} \in Z_n\}$ means the number of reference options used at update step n . For simplicity, we assume this number is always 1.

Assumption A.9. $\sum_{\bar{o}} I\{\bar{o} \in Z_n\} = 1$ for all discrete $n \geq 0$.

The update rules of General Intra-option Differential Q are

$$Q_{n+1}(s, o) \doteq Q_n(s, o) + \alpha_{\nu(n, s, o)} \sum_{\bar{o}} \rho_n(s, o, \bar{o}) \delta_n(s, o, \bar{o}) I\{(s, o) \in Y_n\} I\{\bar{o} \in Z_n\}, \quad \forall s \in \mathcal{S}, \text{ and } o \in \mathcal{O} \quad (\text{A.16})$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta \sum_{s, o} \alpha_{\nu(n, s, o)} \sum_{\bar{o}} \rho_n(s, o, \bar{o}) \delta_n(s, o, \bar{o}) I\{(s, o) \in Y_n\} I\{\bar{o} \in Z_n\}, \quad (\text{A.17})$$

where $\rho_n(s, o, \bar{o}) \doteq \pi(A_n(s, \bar{o}) \mid s, o) / \pi(A_n(s, \bar{o}) \mid s, \bar{o})$ and

$$\begin{aligned} \delta_n(s, o, \bar{o}) &\doteq R_n(s, A_n(s, \bar{o})) - \bar{R}_n + \beta(S'_n(s, A_n(s, \bar{o})), o) \max_{o'} Q_n(S'_n(s, A_n(s, \bar{o})), o') \\ &\quad + (1 - \beta(S'_n(s, A_n(s, \bar{o})), o)) Q_n(S'_n(s, A_n(s, \bar{o})), o) - Q_n(s, o) \end{aligned} \quad (\text{A.18})$$

is the TD error.

Here $\alpha_{\nu(n, s, o)}$ is the stepsize at step n for state-option-option triple (s, o) . The quantity $\alpha_{\nu(n, s, o)}$ depends on the sequence $\{\alpha_n\}$, which is an algorithmic design choice, and also depends on the visitation of state-option pairs $\nu(n, s, o)$. To obtain the stepsize, the algorithm could maintain a $|\mathcal{S} \times \mathcal{O}|$ -size table counting the number of visitations to each state-option pair, which is exactly $\nu(\cdot, \cdot, \cdot)$. Then the stepsize $\alpha_{\nu(n, s, o)}$ can be obtained as long as the sequence $\{\alpha_n\}$ is specified.

Now we show Intra-option Differential Q-learning and Intra-option Differential Q-planning are special cases of General Intra-option Differential Q. Consider a sequence of real experience $\dots, S_t, O_t, A_t, R_{t+1}, S_{t+1}, \dots$. By choosing step $n =$ time step t ,

$$\begin{aligned} Y_n(s, o) &= 1, \text{ if } s = S_t \\ Y_n(s, o) &= 0 \text{ otherwise,} \\ Z_n(\bar{o}) &= 1, \text{ if } \bar{o} = O_t \\ Z_n(\bar{o}) &= 0 \text{ otherwise,} \end{aligned}$$

and $A_n(S_t, O_t) = A_t$, $S'_n(S_t, A_n(S_t, O_t)) = S_{t+1}$, $R_n(S_t, A_n(S_t, O_t)) = R_{t+1}$, update rules (A.16), (A.17), and (A.18) become

$$\begin{aligned} Q_{t+1}(S_t, o) &\doteq Q_t(S_t, o) + \alpha_{\nu(t, S_t, o)} \rho_t(o) \delta_t(o), \forall o \in \mathcal{O}, \text{ and } Q_{t+1}(s, o) \doteq Q_t(s, o), \forall o \in \mathcal{O} \text{ and } \forall s \neq S_t, \\ \bar{R}_{t+1} &\doteq \bar{R}_t + \eta \sum_o \alpha_{\nu(t, S_t, o)} \rho_t(o) \delta_t(o), \\ \delta_t(o) &\doteq R_{t+1} - \bar{R}_t + \beta(S_{t+1}, o) \max_{o'} Q_t(S_{t+1}, o') + (1 - \beta(S_{t+1}, o)) Q_t(S_{t+1}, o) - Q_t(S_t, o), \end{aligned}$$

where $\rho_t(o) \doteq \pi(A_t \mid S_t, o) / \pi(A_t \mid S_t, O_t)$. The above equations are Intra-option Differential Q-learning's update rules (Equations 13, 14, 15) with stepsize at time t being $\alpha_{\nu(t, S_t, o)}$ for each option o .

If we consider a sequence of simulated experience $\dots, \tilde{S}_n, \tilde{O}_n, \tilde{A}_n, \tilde{R}_n, \tilde{S}'_n, \dots$, by choosing step $n =$ planning step n ,

$$\begin{aligned} Y_n(s, o) &= 1, \text{ if } s = \tilde{S}_n \\ Y_n(s, o) &= 0 \text{ otherwise,} \\ Z_n(\bar{o}) &= 1, \text{ if } \bar{o} = \tilde{O}_n \\ Z_n(\bar{o}) &= 0 \text{ otherwise,} \end{aligned}$$

795 and $A_n(\tilde{S}_n, \tilde{O}_n) = \tilde{A}_n$, $S'_n(\tilde{S}_n, A_n(\tilde{S}_n, \tilde{O}_n)) = \tilde{S}'_n$, $R_n(\tilde{S}_n, A_n(\tilde{S}_n, \tilde{O}_n)) = \tilde{R}_n$, update rules
 796 (A.16), (A.17), and (A.18) become

$$\begin{aligned} Q_{n+1}(\tilde{S}_n, o) &\doteq Q_n(\tilde{S}_n, o) + \alpha_{\nu(n, \tilde{S}_n, o)} \rho_n(o) \delta_n(o), \forall o \in \mathcal{O}, \text{ and } Q_{n+1}(s, o) \doteq Q_n(s, o), \forall s \neq \tilde{S}_n, \forall o \in \mathcal{O} \\ \bar{R}_{n+1} &\doteq \bar{R}_n + \eta \sum_o \alpha_{\nu(n, \tilde{S}_n, o)} \rho_n(o) \delta_n(o), \\ \delta_n(o) &\doteq \tilde{R}_n - \bar{R}_n + \beta(\tilde{S}'_n, o) \max_{o'} Q_n(\tilde{S}'_n, o') + (1 - \beta(\tilde{S}'_n, o)) Q_n(\tilde{S}'_n, o) - Q_n(\tilde{S}_n, o), \end{aligned}$$

797 where $\rho_n(o) \doteq \pi(A_n | S_n, o) / \pi(A_n | S_n, O_n)$. The above equations are Intra-option Differential
 798 Q-planning's update rules (Equations 13, 14, 15) with stepsize at planning step n being $\alpha_{\nu(n, S_n, o)}$
 799 for each option o .

800 Finally, note that for both Intra-option Differential Q-learning and Q-planning algorithms, because
 801 for each time step t or update step n , there is only one option which is actually chosen to generate
 802 data, Assumption A.9 is satisfied.

803 **Theorem A.3.** *Under Assumptions 1, A.5, A.6, A.7, A.9, General Intra-option Differential Q (Equa-*
 804 *tions A.16-A.18) converges, almost surely, Q_n to q satisfying both (10)-(11) and*

$$\eta(\sum q - \sum Q_0) = r_* - \bar{R}_0, \quad (\text{A.19})$$

805 \bar{R}_n to r_* , and $r(\mu_n)$ to r_* where μ_n is a greedy policy w.r.t. Q_n .

806 *Proof.* At each step, the increment to \bar{R}_n is η times the increment to Q_n and $\sum Q_n$. Therefore, the
 807 cumulative increment can be written as:

$$\begin{aligned} \bar{R}_n - \bar{R}_0 &= \eta \sum_{i=0}^{n-1} \sum_{s,o} \alpha_{\nu(i, s, o)} \sum_{\bar{o}} \rho_i(s, o, \bar{o}) \delta_i(s, o, \bar{o}) I\{(s, o) \in Y_i\} I\{\bar{o} \in Z_i\} \\ &= \eta \left(\sum Q_n - \sum Q_0 \right) \\ \implies \bar{R}_n &= \eta \sum Q_n - \eta \sum Q_0 + \bar{R}_0 = \eta \sum Q_n - c, \quad (\text{A.20}) \\ \text{where } c &\doteq \eta \sum Q_0 - \bar{R}_0. \end{aligned}$$

808 Now substituting \bar{R}_n in (A.16) with (A.20), we have $\forall s \in \mathcal{S}, o \in \mathcal{O}$:

$$\begin{aligned} Q_{n+1}(s, o) &= Q_n(s, o) + \alpha_{\nu(n, s, o)} \sum_{\bar{o}} \frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} \\ &\quad \left(R_n(s, A_n(s, \bar{o})) - \eta \sum Q_n + c + \beta(S'_n(s, A_n(s, \bar{o})), o) \max_{o'} Q_n(S'_n(s, A_n(s, \bar{o})), o') \right. \\ &\quad \left. + (1 - \beta(S'_n(s, A_n(s, \bar{o})), o)) Q_n(S'_n(s, A_n(s, \bar{o})), o) - Q_n(s, o) \right) \\ &\quad I\{(s, o) \in Y_n\} I\{\bar{o} \in Z_n\}. \quad (\text{A.21}) \end{aligned}$$

809 We now show that (A.21) is a special case of (A.1). To see this point, let $i = (s, o)$,

$$\begin{aligned} R_n(i) &= \sum_{\bar{o}} \frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} I\{\bar{o} \in Z_n\} (R_n(s, A_n(s, \bar{o})) + c), \\ F_n(Q_n)(i) &= \sum_{\bar{o}} \frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} I\{\bar{o} \in Z_n\} \eta \sum Q_n, \\ G_n(Q_n)(i) &= \sum_{\bar{o}} \frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} I\{\bar{o} \in Z_n\} \left(\beta(S'_n(s, A_n(s, \bar{o})), o) \max_{o'} Q_n(S'_n(s, A_n(s, \bar{o})), o') \right. \\ &\quad \left. + (1 - \beta(S'_n(s, A_n(s, \bar{o})), o)) Q_n(S'_n(s, A_n(s, \bar{o})), o) - Q_n(s, o) \right), \\ \epsilon_n(i) &= 0. \end{aligned}$$

810 Then we have:

$$\begin{aligned}
r(i) &= \mathbb{E}[R_n(i)] \\
&= \mathbb{E} \left[\sum_{\bar{o}} \frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} I\{\bar{o} \in Z_n\} (R_n(s, A_n(s, \bar{o})) + c) \right] \\
&= \sum_{\bar{o}} \mathbb{E} \left[\frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} I\{\bar{o} \in Z_n\} (R_n(s, A_n(s, \bar{o})) + c) \right] \\
&= \sum_{\bar{o}} I\{\bar{o} \in Z_n\} \sum_a \pi(a | s, o) \mathbb{E}[R_n(s, a) + c] \\
&= \sum_a \pi(a | s, o) \sum_{r, s'} p(r, s' | s, a) (r + c), \quad \text{By Assumption A.9,} \\
f(q) &= \mathbb{E}[F(q)(i)] = \eta \sum q, \\
g(q)(i) &= \mathbb{E}[G_n(q)(i)] \\
&= \mathbb{E} \left[\sum_{\bar{o}} \frac{\pi(A_n(s, \bar{o}) | s, o)}{\pi(A_n(s, \bar{o}) | s, \bar{o})} I\{\bar{o} \in Z_n\} (\beta(S'_n(s, A_n(s, \bar{o})), o) \max_{o'} q(S'_n(s, A_n(s, \bar{o})), o') \right. \\
&\quad \left. + (1 - \beta(S'_n(s, A_n(s, \bar{o})), o)) q(S'_n(s, A_n(s, \bar{o})), o) - q(s, o)) \right] \\
&= \sum_{\bar{o}} I\{\bar{o} \in Z_n\} \sum_a \pi(a | s, o) \\
&\quad \mathbb{E}[(\beta(S'_n(s, a), o) \max_{o'} q(S'_n(s, a), o') + (1 - \beta(S'_n(s, a), o)) q(S'_n(s, a), o) - q(s, o))] \\
&= \sum_a \pi(a | s, o) \sum_{s', r} p(s', r | s, a) (\beta(s', o) \max_{o'} q(s', o') + (1 - \beta(s', o)) q(s', o) - q(s, o)),
\end{aligned}$$

811 for any $i \in \mathcal{I}$.

812 We now verify the assumptions of Theorem A.1 for Intra-option General Differential
813 Q. Assumption A.1 can be verified for $g(q)(s, o) = \sum_a \pi(a | s, o) \sum_{s', r} p(s', r |$
814 $s, a) (\beta(s', o) \max_{o'} q(s', o') + (1 - \beta(s', o)) q(s', o))$ easily. Assumption A.2 is satisfied for
815 $f(q) = \eta \sum q$. Assumption A.3 satisfies because the MDP is finite. Assumption A.4 is satis-
816 fied for $\epsilon_n = 0$. Assumption A.5-A.7 are satisfied due to assumptions of the theorem being proved.
817 Assumption A.8 is satisfied because

$$\begin{aligned}
&r(i) - \bar{r} + g(q)(i) - q(i) \\
&= \sum_a \pi(a | s, o) \sum_{s', r} p(s', r | s, a) (r - \bar{r} + \beta(s', o) \max_{o'} q(s', o') + (1 - \beta(s', o)) q(s', o)).
\end{aligned}$$

818 By Theorem 2, we know that if the above equation equals to 0, then under Assumption 1, $\bar{r} = r_* + c$
819 is the unique solution and the solutions for q form a set $q = q_* + ke$ for all $k \in \mathbb{R}$.

820 Therefore Theorem A.1 applies and we conclude that Q_n converges to a point satisfying $\eta \sum q =$
821 $r_* + c = r_* + \eta \sum Q_0 - \bar{R}_0$ and $\bar{R}_n = \eta \sum Q_n - c$ to $\eta \sum q - c = r_* + c - c = r_*$. Finally, by
822 Lemma A.6, we have $r(\mu_n) \rightarrow r_*$.

823 □

824 A.5 Theorem 4

825 *Proof.* We will show that there exists a unique solution for (17). Results for (18) and (19) can
826 be shown in a similar way. To show that (17) has a unique solution, we apply a generalized
827 version of the Banach fixed point theorem (see, e.g., Theorem 2.4 by Almezal, Ansari, and Khamsi
828 2014). Once the unique existence of the solution is shown, we easily verify that m^p is the unique
829 solution by showing that it is one solution to (17) as follows. With a little abuse of notation, let

830 $\hat{p}(s', r \mid s, o) \doteq \sum_{r,l} \hat{p}(x, r, l \mid s, o)$, we have

$$\begin{aligned}
m^p(x|s, o) &= \sum_{r,l} \hat{p}(x, r, l|s, o) \\
&= \sum_{l=1}^{\infty} \hat{p}(x, l|s, o) = \sum_a \pi(a|s, o) \sum_r p(s', r|s, a) \beta(s', o) \mathbb{I}(x = s') + \sum_{l=2}^{\infty} \hat{p}(x, l|s, o) \\
&= \sum_a \pi(a|s, o) \sum_r p(s', r|s, a) (\beta(s', o) \mathbb{I}(x = s') + (1 - \beta(s', o)) \sum_{l=1}^{\infty} \hat{p}(x, l|s', o)) \\
&= \sum_a \pi(a|s, o) \sum_r p(s', r|s, a) (\beta(s', o) \mathbb{I}(x = s') + (1 - \beta(s', o)) m^p(x|s', o)).
\end{aligned}$$

831 To apply the generalized version of the Banach fixed point theorem to show the unique
832 existence of the solution, we first define operator $T : \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{O}|} \rightarrow \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{O}|}$
833 such that for any $m \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{O}|}$ and any $x, s \in \mathcal{S}, o \in \mathcal{O}$, $Tm(x \mid s, o) \doteq$
834 $\sum_a \pi(a|s, o) \sum_{s',r} p(s', r|s, a) (\beta(s', o) \mathbb{I}(x = s') + (1 - \beta(s', o)) m(x|s', o))$. We further define
835 $T^n m \doteq T(T^{n-1} m)$ for any $n \geq 2$ and any $m \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{O}|}$. The generalized Banach fixed point
836 theorem shows that if T^n is a contraction mapping for any integer $n \geq 1$ (this is called a n -stage
837 contraction), then $Tm = m$ has a unique fixed point. The unique fixed point immediately leads to
838 the existence of the unique solution of (17). The existence of the unique solution and that m^p is a
839 solution imply that m^p is the unique solution.

840 The only work left is to verify the following contraction property:

$$\|T^{|\mathcal{S}|} m_1 - T^{|\mathcal{S}|} m_2\|_{\infty} \leq \gamma \|m_1 - m_2\|_{\infty}, \quad (\text{A.22})$$

841 where m_1 and m_2 are arbitrary members in $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{O}|}$, and $\gamma < 1$ is some constant.

842 Consider the difference between $T^{|\mathcal{S}|} m_1$ and $T^{|\mathcal{S}|} m_2$ for arbitrary $m_1, m_2 \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{O}|}$. For any
843 $x, s \in \mathcal{S}, o \in \mathcal{O}$, we have

$$\begin{aligned}
&T^{|\mathcal{S}|} m_1(x \mid s, o) - T^{|\mathcal{S}|} m_2(x \mid s, o) \\
&= \sum_a \pi(a \mid s, o) \sum_{s',r} p(s', r|s, a) (1 - \beta(s', o)) (T^{|\mathcal{S}|-1} m_1(x \mid s', o) - T^{|\mathcal{S}|-1} m_2(x \mid s', o)) \\
&= \sum_{s_1} \Pr(S_{t+1} = s_1 \mid S_t = s, O_t = o) (1 - \beta(s_1, o)) (T^{|\mathcal{S}|-1} m_1(x \mid s_1, o) - T^{|\mathcal{S}|-1} m_2(x \mid s_1, o)) \\
&= \sum_{s_1} \Pr(S_{t+1} = s_1 \mid S_t = s, O_t = o) (1 - \beta(s_1, o)) \sum_{s_2} \Pr(S_{t+2} = s_2 \mid S_{t+1} = s_1, O_{t+1} = o) (1 - \beta(s_2, o)) \\
&\quad (T^{|\mathcal{S}|-2} m_1(x \mid s_2, o) - T^{|\mathcal{S}|-2} m_2(x \mid s_2, o)) \\
&\vdots \\
&= \sum_{s_1, \dots, s_{|\mathcal{S}|}} \Pr(S_{t+1} = s_1, \dots, S_{t+|\mathcal{S}|} = s_{|\mathcal{S}|} \mid S_t = s, O_t = o) \prod_{i=1}^{|\mathcal{S}|} (1 - \beta(s_i, o)) (m_1(x \mid s_{|\mathcal{S}|}, o) - m_2(x \mid s_{|\mathcal{S}|}, o)) \\
&\leq \sum_{s_1, \dots, s_{|\mathcal{S}|}} \Pr(S_{t+1} = s_1, \dots, S_{t+|\mathcal{S}|} = s_{|\mathcal{S}|} \mid S_t = s, O_t = o) \prod_{i=1}^{|\mathcal{S}|} (1 - \beta(s_i, o)) \|m_1 - m_2\|_{\infty}.
\end{aligned}$$

844 Here $\tilde{p}(s, o) \doteq \sum_{s_1, \dots, s_{|\mathcal{S}|}} \Pr(S_{t+1} = s_1, \dots, S_{t+|\mathcal{S}|} = s_{|\mathcal{S}|} \mid S_t = s, O_t = o) \prod_{i=1}^{|\mathcal{S}|} (1 - \beta(s_i, o))$
845 is the probability of executing option o for $|\mathcal{S}|$ steps starting from s without termination. If $\tilde{p}(s, o) = 0$,
846 then option o will surely terminate within the first $|\mathcal{S}|$ steps and if $\tilde{p}(s, o) = 1$, then option o will
847 surely not terminate within the first $|\mathcal{S}|$ steps.

848 If option o would surely not terminate within the first $|\mathcal{S}|$ steps ($\tilde{p}(s, o) = 1$), then it would surely not
849 terminate forever. This is because there are only $|\mathcal{S}|$ number of states, and thus an option could visit
850 all states that are possible to be visited by the option within the first $|\mathcal{S}|$ steps. $\tilde{p}(s, o) = 1$ means that

option o has a zero probability of terminating in all states that are possible to be visited by option o . This non-termination of a state-option pair implies that the expected option length is infinite, which is contradict to our assumption of finite expected option lengths (Section 2). Therefore $\tilde{p}(s, o) = 1$ is not allowed by our assumption and thus $\tilde{p}(s, o) < 1$. So there must exist some $\gamma(s, o) < 1$ such that $\tilde{p}(s, o) \leq \gamma(s, o)$. With $\gamma \doteq \max_{s,o} \gamma(s, o)$, we obtain (A.22). \square

A.6 Theorem 5

We first provide a formal statement of Theorem 5. The formal theorem statement needs stepsizes to be specific for each state-option pair. We rewrite (20–22) to incorporate such stepsizes:

$$M_{t+1}^p(x | S_t, o) \doteq M_t^p(x | S_t, o) + \alpha_t(S_t, o)\rho_t(o)\left(\beta(S_{t+1}, o)\mathbb{I}(S_{t+1} = x) + (1 - \beta(S_{t+1}, o))M_t^p(x | S_{t+1}, o) - M_t^p(x | S_t, o)\right), \quad \forall x \in \mathcal{S}, \quad (\text{A.23})$$

$$M_{t+1}^r(S_t, o) \doteq M_t^r(S_t, o) + \alpha_t(S_t, o)\rho_t(o)\left(R_{t+1} + (1 - \beta(S_{t+1}, o))M_t^r(S_{t+1}, o) - M_t^r(S_t, o)\right) \quad (\text{A.24})$$

$$M_{t+1}^l(S_t, o) \doteq M_t^l(S_t, o) + \alpha_t(S_t, o)\rho_t(o)\left(1 + (1 - \beta(S_{t+1}, o))M_t^l(S_{t+1}, o) - M_t^l(S_t, o)\right). \quad (\text{A.25})$$

Theorem A.4 (Convergence of the intra-option model learning algorithm, formal). *If $0 \leq \alpha_t(s, o) \leq 1$, $\sum_t \alpha_t(s, o) = \infty$ and $\sum_t \alpha_t^2(s, o) < \infty$, and $\alpha_t(s, o) = 0$ unless $s = S_t$, then the intra-option model-learning algorithm (A.23–A.25) converges almost surely, M_t^p to m^p , M_t^r to m^r , and M_t^l to m^l .*

Here the assumptions on α_t guarantee that each state-option pair is updated for an infinite number of times. Because the three update rules are independent, we only show convergence of the first update rule; the other two can be shown in the same way.

Proof. We apply a slight generalization of Theorem 3 by Tsitsiklis (1994) to show the above theorem. The generalization replaces Assumption 5 (an assumption for Theorem 3) by:

Assumption A.10. *There exists a vector $x^* \in \mathbb{R}^n$, a positive vector v , a positive integer m , and a scalar $\beta \in [0, 1]$, such that*

$$\|F^m(x) - x^*\|_v \leq \beta \|x - x^*\|_v, \quad \forall x \in \mathbb{R}^n.$$

That is, we replace the one-stage contraction assumption by a m -stage contraction assumption. The proof of Tsitsiklis’ Theorem 3 also applies to its generalized form and is thus omitted here.

Notice that our update rule (A.23) is a special case of the general update rule considered by Theorem 3 (equations 1-3), and is thus a special case of its generalized version. Therefore we only need to verify the above m -stage contraction assumption, as well as Assumption 1, 2, and 3 required by Theorem 3. According to the proof in Appendix A.5, the operator T associated with the update rule (20) is a $|\mathcal{S}|$ -stage contraction (and thus is a $|\mathcal{S}|$ -stage pseudo-contraction). Other assumptions (Assumptions 1, 2, 3) required by Theorem 3 are also satisfied given our step-size, and finite MDP assumptions. \square

A.7 Theorem 6

Proof. We first show that

$$\begin{aligned} & \sum_{o'} \mu'(o' | s) \sum_{s', r, l} \hat{p}(s', r, l | s, o')(r - lr(\mu) + v_\mu(s')) \\ & \geq \sum_o \mu(o | s) \sum_{s', r, l} \hat{p}(s', r, l | s, o)(r - lr(\mu) + v_\mu(s')) = v_\mu(s). \end{aligned} \quad (\text{A.26})$$

Note that for all s, o and its corresponding o' , $\mu(o | s) = \mu'(o' | s)$. In order to show (A.26), we show $\sum_{s', r, l} \hat{p}(s', r, l | s, o')(r - lr(\mu) + v_\mu(s')) \geq \sum_{s', r, l} \hat{p}(s', r, l | s, o)(r - lr(\mu) + v_\mu(s'))$ for

all s, o and corresponding o' .

$$\begin{aligned}
& \sum_{s', r, l} \hat{p}(s', r, l | s, o')(r - lr(\mu) + v_\mu(s')) \\
&= \mathbb{E}[\hat{R}_n - \hat{L}_n r(\mu) + v_\mu(\hat{S}_{n+1}) | S_n = s, O_n = o'] \\
&= \mathbb{E}[\hat{R}_n - \hat{L}_n r(\mu) + v_\mu(\hat{S}_{n+1}) | S_n = s, O_n = o', \text{Not encountering an interruption}] \\
&+ \mathbb{E}[\hat{R}_n - \hat{L}_n r(\mu) + v_\mu(\hat{S}_{n+1}) | S_n = s, O_n = o', \text{Encountering an interruption}] \\
&\geq \mathbb{E}[\hat{R}_n - \hat{L}_n r(\mu) + v_\mu(\hat{S}_{n+1}) | S_n = s, O_n = o', \text{Not encountering an interruption}] \\
&+ \mathbb{E}[\beta(s')(\hat{R}_n - \hat{L}_n r(\mu) + v_\mu(\hat{S}_{n+1})) + (1 - \beta(s'))(\hat{R}_n - \hat{L}_n r(\mu) + q_\mu(\hat{S}_{n+1}, o)) \\
&| S_n = s, O_n = o', \text{Encountering an interruption}] \\
&= \sum_{s', r, l} \hat{p}(s', r, l | s, o)(r - lr(\mu) + v_\mu(s')).
\end{aligned}$$

The above inequality holds because \hat{S}_{n+1} is the state where termination happens and thus $q_\mu(\hat{S}_{n+1}, o) \leq v_\mu(\hat{S}_{n+1})$. The last equality holds because $\mathbb{E}[\beta(s')(\hat{R}_n - \hat{L}_n r(\mu) + v_\mu(\hat{S}_{n+1})) + (1 - \beta(s'))(\hat{R}_n - \hat{L}_n r(\mu) + q_\mu(\hat{S}_{n+1}, o)) | S_n = s, O_n = o', \text{Encountering an interruption}]$ is the expected differential return when the agent could interrupt its old option but chooses to stick on the old option. (A.26) is shown.

Now write the l.h.s. of (A.26) in the matrix form

$$\sum_{o'} \mu'(o' | s) \sum_{s', r, l} \hat{p}(s', r, l | s, o')(r - lr(\mu) + v_\mu(s')) = r_{\mu'}(s) - l_{\mu'}(s)r(\mu) + (P_{\mu'} v_\mu)(s),$$

where $r_{\mu'}(s) \doteq \sum_{o'} \mu'(o' | s) \sum_{s', r, l} \hat{p}(s', r, l | s, o')r$ is the expected one option-transition reward, $l_{\mu'}(s) \doteq \sum_{o'} \mu'(o' | s) \sum_{s', r, l} \hat{p}(s', r, l | s, o')l$ is the expected one option-transition length, and $P_{\mu'}(s, s') \doteq \sum_{o'} \mu'(o' | s) \sum_{r, l} \hat{p}(s', r, l | s, o')$ is the probability of terminating at s' .

Combined with the r.h.s. of (A.26), we have

$$r_{\mu'}(s) - l_{\mu'}(s)r(\mu) + (P_{\mu'} v_\mu)(s) \geq v_\mu(s).$$

Iterating the above inequality for $K - 1$ times, we have

$$\begin{aligned}
& \sum_{k=0}^{K-1} (P_{\mu'}^k r_{\mu'}(s) - P_{\mu'}^k l_{\mu'}(s)r(\mu)) + P_{\mu'}^K v_\mu(s) \geq v_\mu(s) \\
& \sum_{k=0}^{K-1} (P_{\mu'}^k r_{\mu'}(s) - P_{\mu'}^k l_{\mu'}(s)r(\mu)) \geq v_\mu(s) - P_{\mu'}^K v_\mu(s).
\end{aligned}$$

Divide both sides by $\sum_{k=0}^{K-1} P_{\mu'}^k l_{\mu'}(s)$ and take $K \rightarrow \infty$:

$$\lim_{K \rightarrow \infty} \frac{1}{\sum_{k=0}^{K-1} P_{\mu'}^k l_{\mu'}(s)} \sum_{k=0}^{K-1} (P_{\mu'}^k r_{\mu'}(s) - P_{\mu'}^k l_{\mu'}(s)r(\mu)) \geq \lim_{K \rightarrow \infty} \frac{1}{\sum_{k=0}^{K-1} P_{\mu'}^k l_{\mu'}(s)} (v_\mu(s) - P_{\mu'}^K v_\mu(s)).$$

For the l.h.s.:

$$\lim_{K \rightarrow \infty} \frac{1}{\sum_{k=0}^{K-1} P_{\mu'}^k l_{\mu'}(s)} \sum_{k=0}^{K-1} (P_{\mu'}^k r_{\mu'}(s) - P_{\mu'}^k l_{\mu'}(s)r(\mu)) = \lim_{K \rightarrow \infty} \frac{\sum_{k=0}^{K-1} P_{\mu'}^k r_{\mu'}(s)}{\sum_{k=0}^{K-1} P_{\mu'}^k l_{\mu'}(s)} - r(\mu) = r(\mu') - r(\mu).$$

For the r.h.s.:

$$\lim_{K \rightarrow \infty} \frac{1}{\sum_{k=0}^{K-1} P_{\mu'}^k l_{\mu'}(s)} (v_\mu(s) - P_{\mu'}^K v_\mu(s)) = 0.$$

Therefore $r(\mu') - r(\mu) \geq 0$.

Finally, note that a strict inequality holds if the probability of interruption when following policy μ' is non-zero. \square

901 B Additional Empirical Results

902 B.1 Inter-option Learning

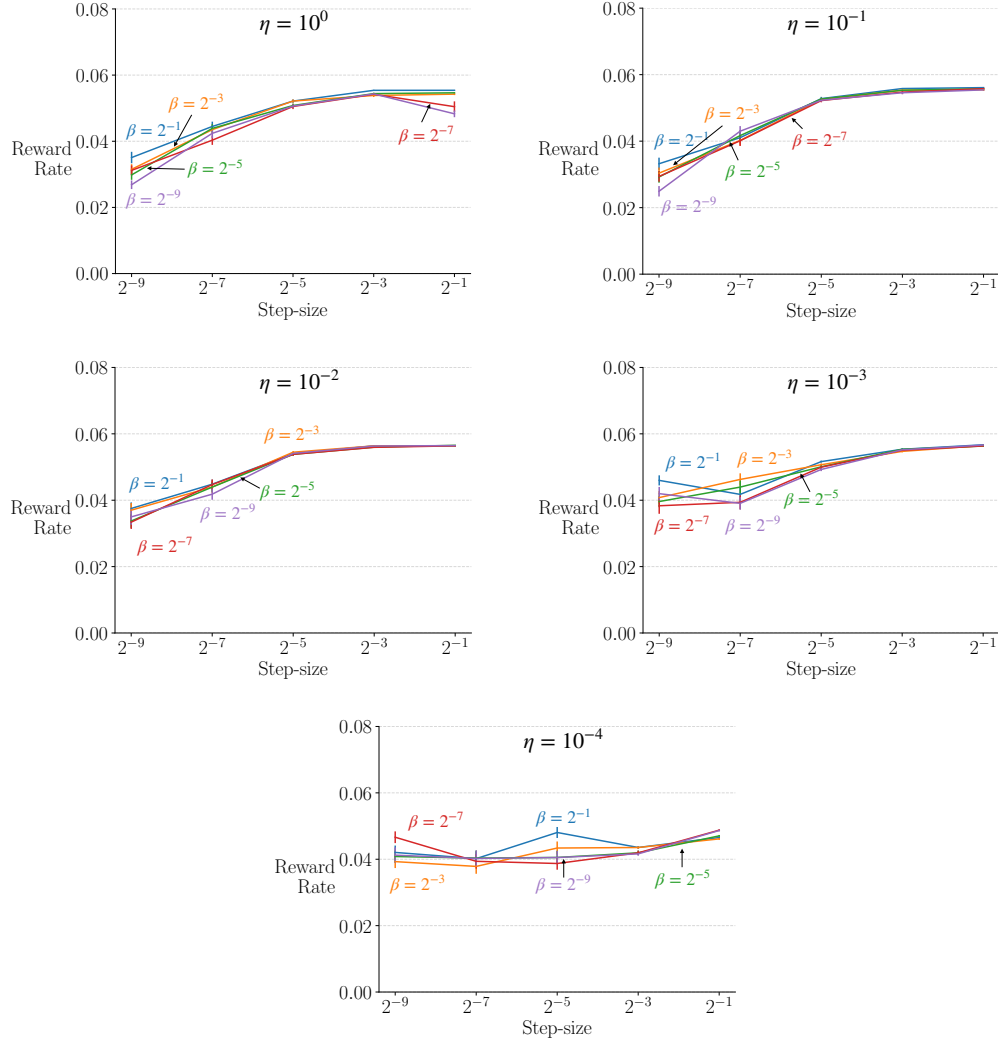


Figure B.1: Plots showing a parameter study for inter-option Differential Q-learning and the set of options $\mathcal{O} = \mathcal{H} + \mathcal{A}$ in the continuing Four-Room domain when the goal was to go to G1. Same experimental setups are used as what was described in Section 3. The x-axis denotes step size α ; the y-axis denotes the rate of the rewards averaged over all 200,000 steps of training, reflecting the rate of learning. The error bars denote one standard error. The algorithm's rate of learning varied little over a broad range of its parameters α , β and η . Small standard error bars show that the algorithm's performance varied little over multiple runs.

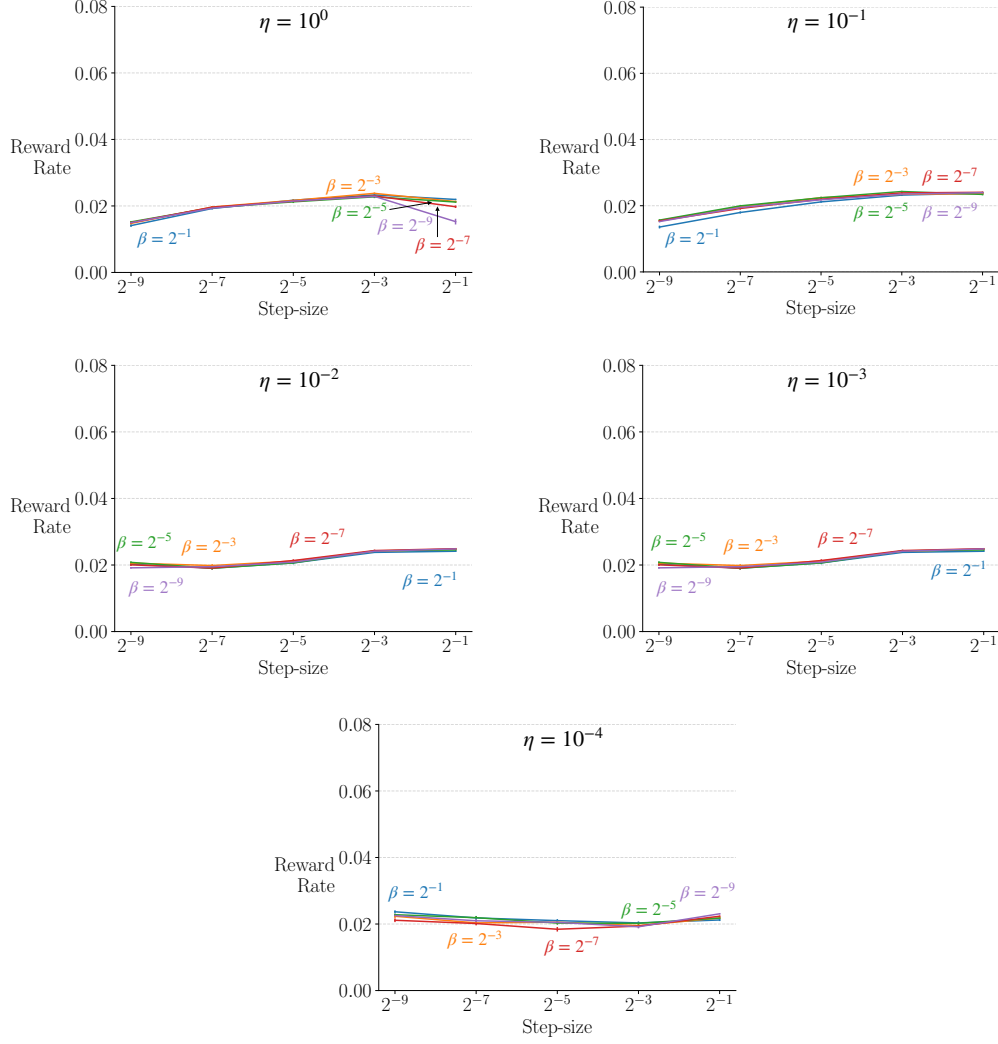


Figure B.2: Plots showing a parameter study for inter-option Differential Q-learning and the set of options $\mathcal{O} = \mathcal{H}$ in the continuing Four-Room domain when the goal was to go to G1. The experimental setting and the plot axes are the same as mentioned in Figure B.1. Compared with Figure B.1, it can be seen that the algorithm's rate of learning with $\mathcal{O} = \mathcal{H}$ was worse than it with $\mathcal{O} = \mathcal{H} + \mathcal{A}$. This is because there is no hallway option from \mathcal{H} can takes the agent to G1. The algorithm's rate of learning varied little over a broad range of its parameters α, β and η , and also varied little over multiple runs.

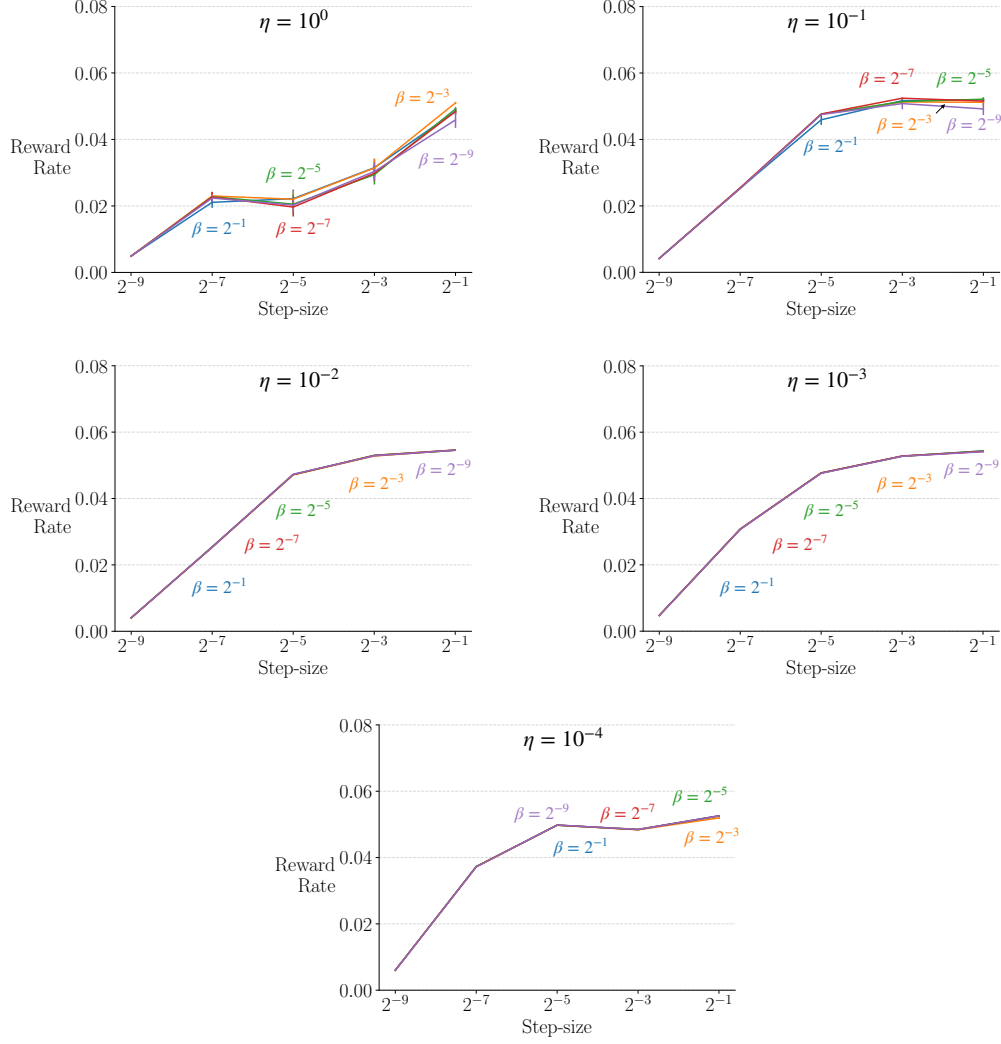


Figure B.3: Plots showing a parameter study for inter-option Differential Q-learning and the set of options $\mathcal{O} = \mathcal{A}$ in the continuing Four-Room domain when the goal was to go to G1. Note that with options being primitive actions, the algorithm becomes exactly the same as Differential Q-learning by Wan et al. (2021). The experimental setting and the plot axes are the same as mentioned in Figure B.1. Compared with Figure B.1, it can be seen that the algorithm’s rate of learning with $\mathcal{O} = \mathcal{A}$ was worse than it with $\mathcal{O} = \mathcal{H} + \mathcal{A}$, particularly for small α . The algorithm’s rate of learning did not vary too much over a broad range of its parameters β and η , and also varied little over multiple runs. The algorithm’s performance is more sensitive to the choice of α .

905 B.2 Intra-option Q-learning

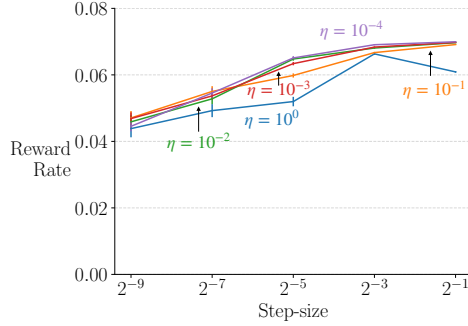


Figure B.4: Plots showing a parameter study for intra-option Differential Q-learning with the set of options $\mathcal{O} = \mathcal{H}$ in the continuing Four-Room domain when the goal was to go to G2. The algorithm used a behavior policy consisting only of primitive actions. The hallway options were never executed.. The experimental setting and the plot axes are the same as mentioned in Section 4. The algorithm’s rate of learning varied little over a broad range of its parameters α and η , and also varied little over multiple runs.

906 B.3 Interruption

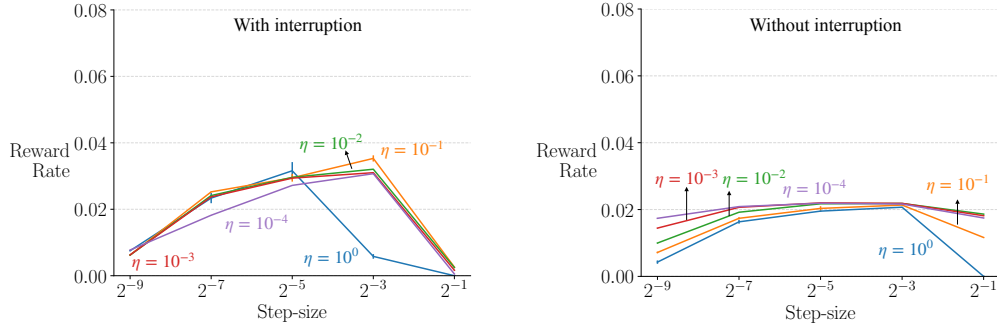


Figure B.5: Plots showing parameter studies for intra-option Differential Q-learning with and without interruption in the continuing Four-Room domain when the goal was to go to G3. The algorithm used the set of hallway options $\mathcal{O} = \mathcal{H}$. The experimental setting and the plot axes are the same as mentioned in Section 6. The algorithm’s rate of learning with interruption was higher than it without interruption for medium sized choices of α . When a large or small α was used, interruption produced a worse rate of learning. The algorithm’s rate of learning varied not too much over a broad range of its parameters η and varied little over multiple runs, regardless of interruption. The algorithm’s rate of learning was more sensitive to α when interruption is used.

907 C Additional Discussion

908 C.1 Two Failed Attempts on Extending Differential Q-learning to an Inter-option Algorithm

909 The authors have tried two other ways of extending Differential Q-learning to an Inter-option
 910 Algorithm (cf. Section 3). While these two ways appeared to work properly at the first glance, they
 911 do not actually. We now show these two approaches and explain why they do not work properly.

912 The first extension uses, for each option, the average-reward rate per-step instead of the total reward
 913 as the reward of the option. In particular, such an extension use update rules (3) and (4), but with TD
 914 error defined as:

$$\delta'_n \doteq \hat{R}_n / \hat{L}_n - \bar{R}_n + \max_o Q_n(\hat{S}_{n+1}, o) - Q_n(\hat{S}_n, \hat{O}_n) \quad (\text{C.1})$$

915 Unfortunately, such an extension can not guarantee convergence to a desired point. Specifically, the
 916 extension, if converges, will converge to a solution of $\mathbb{E}[\delta'_n] = 0$, which is not necessarily a solution
 917 of the Bellman equation $\mathbb{E}[\delta_n] = 0$ (Equation 2).

918 An alternative approach to avoid the instability issue is to shrink the entire update, not the option's
 919 cumulative reward, by the sample length:

$$Q_{n+1}(\hat{S}_n, \hat{O}_n) \doteq Q_n(\hat{S}_n, \hat{O}_n) + \alpha_n \delta_n / \hat{L}_n, \quad (\text{C.2})$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta \alpha_n \delta_n / \hat{L}_n. \quad (\text{C.3})$$

920 Still, the above two updates can not guarantee convergence to the desired values because, again,
 921 $\mathbb{E}[\delta_n / \hat{L}_n] = 0$ does not imply that the Bellman equation $\mathbb{E}[\delta_n] = 0$ is satisfied.

922 C.2 Pseudocodes

Algorithm 1: Inter-option Differential Q-learning

Input: Behavioral policy μ_b 's parameters (e.g., ϵ for ϵ -greedy)

Algorithm parameters: step-size parameters α, η, β

```

1 Initialize  $Q(s, o) \forall s \in \mathcal{S}, o \in \mathcal{O}$ ,  $\bar{R}$  arbitrarily (e.g., to zero);  $L(s, o) \leftarrow 1 \forall s \in \mathcal{S}, o \in \mathcal{O}$ 
2 Obtain initial  $S$ 
3 while still time to train do
4   Initialize  $\hat{L} \leftarrow 0, \hat{R} \leftarrow 0, S_{tmp} \leftarrow S$ 
5    $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$ 
6   do
7     Sample primitive action  $A \sim \pi(\cdot | S, O)$ 
8     Take action  $A$ , observe  $R, S'$ 
9      $\hat{L} \leftarrow \hat{L} + 1$ 
10     $\hat{R} \leftarrow \hat{R} + R$ 
11     $S \leftarrow S'$ 
12  while  $O$  doesn't terminate in  $S'$ 
13   $S \leftarrow S_{tmp}$ 
14   $L(S, O) \leftarrow L(S, O) + \beta(\hat{L} - L(S, O))$ 
15   $\delta \leftarrow \hat{R} - \bar{R} \cdot L(S, O) + \max_o Q(S', o) - Q(S, O)$ 
16   $Q(S, O) \leftarrow Q(S, O) + \alpha \delta / L(S, O)$ 
17   $\bar{R} \leftarrow \bar{R} + \eta \alpha \delta / L(S, O)$ 
18   $S \leftarrow S'$ 
19 end
20 return  $Q$ 

```

Algorithm 2: Inter-option Differential Q-evaluation (learning)

Input: Behavioral policy μ_b , target policy μ

Algorithm parameters: step-size parameters α, η, β

```
1 Initialize  $Q(s, o) \forall s \in \mathcal{S}, o \in \mathcal{O}, \bar{R}$  arbitrarily (e.g., to zero);  $L(s, o) \leftarrow 1 \forall s \in \mathcal{S}, o \in \mathcal{O}$ 
2 Obtain initial  $S$ 
3 while still time to train do
4   Initialize  $\hat{L} \leftarrow 0, \hat{R} \leftarrow 0, S_{tmp} \leftarrow S$ 
5    $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$ 
6   do
7     Sample primitive action  $A \sim \pi(\cdot | S, O)$ 
8     Take action  $A$ , observe  $R, S'$ 
9      $\hat{L} \leftarrow \hat{L} + 1$ 
10     $\hat{R} \leftarrow \hat{R} + R$ 
11     $S \leftarrow S'$ 
12  while  $O$  doesn't terminate in  $S'$ 
13   $S \leftarrow S_{tmp}$ 
14   $L(S, O) \leftarrow L(S, O) + \beta(\hat{L} - L(S, O))$ 
15   $\delta \leftarrow \hat{R} - \bar{R} \cdot L(S, O) + \sum_o \mu(o | S')Q(S', o) - Q(S, O)$ 
16   $Q(S, O) \leftarrow Q(S, O) + \alpha\delta/L(S, O)$ 
17   $\bar{R} \leftarrow \bar{R} + \eta\alpha\delta/L(S, O)$ 
18   $S \leftarrow S'$ 
19 end
20 return  $Q$ 
```

Algorithm 3: Intra-option Differential Q-learning

Input: Behavioral policy μ_b 's parameters (e.g., ϵ for ϵ -greedy)

Algorithm parameters: step-size parameters α, η

```
1 Initialize  $Q(s, o) \forall s \in \mathcal{S}, o \in \mathcal{O}, \bar{R}$  arbitrarily (e.g., to zero)
2 Obtain initial  $S$ 
3 while still time to train do
4    $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$ 
5   do
6     Sample primitive action  $A \sim \pi(\cdot | S, O)$ 
7     Take action  $A$ , observe  $R, S'$ 
8     for all options  $o$  do
9        $\rho \leftarrow \pi(A | S, o)/\pi(A | S, O)$ 
10       $\delta \leftarrow R - \bar{R} + \left( (1 - \beta(S', o))Q(S', o) + \beta(S', o) \max_{o'} Q(S', o') \right) - Q(S, o)$ 
11       $Q(S, o) \leftarrow Q(S, o) + \alpha\rho\delta$ 
12       $\bar{R} \leftarrow \bar{R} + \eta\alpha\rho\delta$ 
13    end
14     $S \leftarrow S'$ 
15  while  $O$  doesn't terminate in  $S$ 
16 end
17 return  $Q$ 
```

Algorithm 4: Intra-option Differential Q-learning with interruption

Input: Behavioral policy μ_b 's parameters (e.g., ϵ for ϵ -greedy)

Algorithm parameters: step-size parameters α, η

```
1 Initialize  $Q(s, o) \forall s \in \mathcal{S}, o \in \mathcal{O}, \bar{R}$  arbitrarily (e.g., to zero)
2 Obtain initial  $S$ 
3  $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$ 
4 while still time to train do
5   if  $O \notin \operatorname{argmax} Q(S, \cdot)$  then
6      $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$ 
7   end
8   Sample primitive action  $A \sim \pi(\cdot | S, O)$ 
9   Take action  $A$ , observe  $R, S'$ 
10  for all options  $o$  do
11     $\rho \leftarrow \pi(A | S, o) / \pi(A | S, O)$ 
12     $\delta \leftarrow R - \bar{R} + \left( (1 - \beta(S', o))Q(S', o) + \beta(S', o) \max_{o'} Q(S', o') \right) - Q(S, o)$ 
13     $Q(S, o) \leftarrow Q(S, o) + \alpha \rho \delta$ 
14     $\bar{R} \leftarrow \bar{R} + \eta \alpha \rho \delta$ 
15  end
16   $S = S'$ 
17 end
18 return  $Q$ 
```

Algorithm 5: Intra-option Differential Q-evaluation (learning)

Input: Behavioral policy μ_b , target policy μ

Algorithm parameters: step-size parameters α, η

```
1 Initialize  $Q(s, o) \forall s \in \mathcal{S}, o \in \mathcal{O}, \bar{R}$  arbitrarily (e.g., to zero)
2 Obtain initial  $S$ 
3 while still time to train do
4    $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$ 
5   do
6     Sample primitive action  $A \sim \pi(\cdot | S, O)$ 
7     Take action  $A$ , observe  $R, S'$ 
8     for all options  $o$  do
9        $\rho \leftarrow \pi(A | S, o) / \pi(A | S, O)$ 
10       $\delta \leftarrow R - \bar{R} + \left( (1 - \beta(S', o))Q(S', o) + \beta(S', o) \sum_{o'} \mu(o' | S')Q(S', o') \right) - Q(S, o)$ 
11       $Q(S, o) \leftarrow Q(S, o) + \alpha \rho \delta$ 
12       $\bar{R} \leftarrow \bar{R} + \eta \alpha \rho \delta$ 
13    end
14     $S \leftarrow S'$ 
15  while  $O$  doesn't terminate in  $S$ 
16 end
17 return  $Q$ 
```

Algorithm 6: Combined Algorithm: Intra-option Model-learning + Inter-option Q-planning

Input: Behavioral policy μ_b 's parameters (e.g., ϵ for ϵ -greedy)

Algorithm parameters: step-size parameters α, β, η ; number of planning steps per time step n

```
1 Initialize  $Q(s, o), P(x | s, o), R(s, o) \forall s, x \in \mathcal{S}, o \in \mathcal{O}, \bar{R}$ , arbitrarily (e.g., to zero);  
    $L(s, o) = 1 \forall s \in \mathcal{S}, o \in \mathcal{O}; T \leftarrow False$   
2 while still time to train do  
3    $S \leftarrow$  current state  
4    $O \leftarrow$  option sampled from  $\mu_b(\cdot | S)$   
5   while  $T$  is False do  
6     Sample primitive action  $A \sim \pi(\cdot | S, O)$   
7     Take action  $A$ , observe  $R', S'$   
8     for all options  $o$  such that  $\pi(A | S, o) > 0$  do  
9        $\rho \leftarrow \pi(A | S, o) / \pi(A | S, O)$   
10      for all states  $x \in \mathcal{S}$  do  
11         $P(x | S, o) \leftarrow P(x | S, o) + \beta \rho (\beta(S', o) \mathbb{I}(S' = x) + (1 - \beta(S', o)) P(x | S', o) - P(x | S, o))$   
12      end  
13       $R(S, o) \leftarrow R(S, o) + \beta \rho (R' + (1 - \beta(S', o)) R(S', o) - R(S, o))$   
14       $L(S, o) \leftarrow L(S, o) + \beta \rho (1 + (1 - \beta(S', o)) L(S', o) - L(S, o))$   
15    end  
16     $T \leftarrow$  indicator of termination sampled from  $\beta(S', O)$   
17    for all of the  $n$  planning steps do  
18       $S \leftarrow$  a random previously observed state  
19       $O \leftarrow$  a random option previously taken in  $S$   
20       $S' \leftarrow$  a sampled state from  $P(\cdot | S, O)$   
21       $\delta \leftarrow R(S, O) - L(S, O) \bar{R} + \max_o Q(S', o) - Q(S, O)$   
22       $Q(S, O) \leftarrow Q(S, O) + \alpha \rho \delta / L(S, O)$   
23       $\bar{R} \leftarrow \bar{R} + \eta \alpha \rho \delta / L(S, O)$   
24    end  
25  end  
26 end  
27 return  $Q$ 
```
