# Burning RED: Unlocking Subtask-Driven Reinforcement Learning and Risk-Awareness in Average-Reward Markov Decision Processes

**Anonymous authors**
Paper under double-blind review

**Keywords:** Average-Reward Reinforcement Learning, Risk-Sensitive Decision-Making, CVaR

## Summary

Average-reward Markov decision processes (MDPs) provide a foundational framework for sequential decision-making under uncertainty. However, average-reward MDPs have remained largely unexplored in reinforcement learning (RL) settings, with the majority of RL-based efforts having been allocated to discounted MDPs. In this work, we study a unique structural property of average-reward MDPs and utilize it to introduce *Reward-Extended Differential* (or *RED*) reinforcement learning: a novel RL framework that can be used to effectively and efficiently solve various learning objectives, or *subtasks*, simultaneously in the average-reward setting. We introduce a family of RED learning algorithms for prediction and control, including proven-convergent algorithms for the tabular case. We then showcase the power of these algorithms by demonstrating how they can be used to learn a policy that optimizes, for the first time, the well-known conditional value-at-risk (CVaR) risk measure in a fully-online manner, *without* the use of an explicit bi-level optimization scheme or an augmented state-space.

## Contribution(s)

1. We provide a general-purpose framework and a corresponding set of prediction/control algorithms for solving an arbitrary number of learning objectives, or *subtasks*, simultaneously in the average-reward setting with only a TD error-based update, including proven-convergent algorithms for the tabular case.
   **Context:** Our work builds on (and can be viewed as a generalization of) Wan et al. (2021), which proposed proven-convergent average-reward RL algorithms that are able to learn and/or optimize the value function and average-reward simultaneously using only the TD error. In particular, the focus in Wan et al. (2021) was on proving the convergence of such algorithms, without exploring the underlying structural properties of the average-reward MDP that made such a process possible to begin with. In this work, we formalize these underlying properties, and utilize them to show that if one modifies, or *extends*, the reward from the MDP with various learning objectives, then these objectives, or *subtasks*, can be solved simultaneously using a modified, or *reward-extended*, version of the TD error.

2. We provide the first RL algorithm that optimizes the well-known conditional value-at-risk (CVaR) risk measure (Rockafellar and Uryasev, 2000) in a fully-online manner *without* the use of an explicit bi-level optimization or an augmented state-space.
   **Context:** Several prior works have looked at CVaR optimization in the discounted setting (e.g. Bäuerle and Ott (2011) and Chow et al. (2015)). However, no prior work has developed an algorithm for CVaR optimization that does not require either an augmented state-space or an explicit bi-level optimization, which can, for example, involve solving multiple MDPs. In the average-reward setting, Xia et al. (2023) proposed a set of algorithms for optimizing the CVaR risk measure, however their methods require the use of an augmented state-space and a sensitivity-based bi-level optimization. By contrast, our work, to the best of our knowledge, is the first to optimize CVaR in an MDP-based setting without the use of an explicit bi-level optimization scheme or an augmented state-space.

# Burning RED: Unlocking Subtask-Driven Reinforcement Learning and Risk-Awareness in Average-Reward Markov Decision Processes

**Anonymous authors**
Paper under double-blind review

## Abstract

Average-reward Markov decision processes (MDPs) provide a foundational framework for sequential decision-making under uncertainty. However, average-reward MDPs have remained largely unexplored in reinforcement learning (RL) settings, with the majority of RL-based efforts having been allocated to discounted MDPs. In this work, we study a unique structural property of average-reward MDPs and utilize it to introduce *Reward-Extended Differential* (or *RED*) reinforcement learning: a novel RL framework that can be used to effectively and efficiently solve various learning objectives, or *subtasks*, simultaneously in the average-reward setting. We introduce a family of RED learning algorithms for prediction and control, including proven-convergent algorithms for the tabular case. We then showcase the power of these algorithms by demonstrating how they can be used to learn a policy that optimizes, for the first time, the well-known conditional value-at-risk (CVaR) risk measure in a fully-online manner, *without* the use of an explicit bi-level optimization scheme or an augmented state-space.

## 1 Introduction

Markov decision processes (MDPs) (Puterman, 1994) are a long-established framework for sequential decision-making under uncertainty. Discounted MDPs, which aim to optimize a potentially-discounted sum of rewards over time, have enjoyed success in recent years when utilizing reinforcement learning (RL) solution methods (Sutton and Barto, 2018) to tackle certain problems of interest in various domains. Despite this success however, these MDP-based methods have yet to be fully embraced in real-world applications due to the various intricacies and implications of real-world operation that often trump the ability of current state-of-the-art methods (Dulac-Arnold et al., 2021). We therefore turn to the less-explored average-reward MDP, which aims to optimize the reward received per time-step, to see how its unique structural properties can be leveraged to tackle challenging problems that have evaded its discounted counterpart.

In particular, we present results that show how the average-reward MDP's unique structural properties can be leveraged to enable a more *subtask-driven* approach to reinforcement learning, where various learning objectives, or *subtasks*, are solved simultaneously (and in a fully-online manner) to help solve a larger, central learning objective. Importantly, we find a compelling case-study in the realm of risk-aware decision-making that illustrates how this subtask-driven approach can alleviate some of the computational challenges and non-trivialities that can arise in the discounted setting.

More formally, we introduce *Reward-Extended Differential* (or *RED*) reinforcement learning: a first-of-its-kind RL framework that makes it possible to solve various subtasks simultaneously in the average-reward setting. At the heart of this framework is the novel concept of the reward-extended temporal-difference (TD) error, an extension of the celebrated TD error (Sutton, 1988), which we derive by leveraging a unique structural property of average-reward MDPs, and utilize to solve

36 various subtasks simultaneously. We first present the RED RL framework in a generalized way, then
37 adopt it to successfully tackle a problem that has exceeded the capabilities of current state-of-the-art
38 methods in risk-aware decision-making: learning a policy that optimizes the well-known conditional
39 value-at-risk (CVaR) risk measure (Rockafellar and Uryasev, 2000) in a fully-online manner *without*
40 the use of an explicit bi-level optimization scheme or an augmented state-space.

41 Our work is organized as follows: In Section 2, we provide a brief overview of related work. In Sec-
42 tion 3, we give an overview of the fundamental concepts related to average-reward RL and CVaR.
43 In Section 4, we motivate the need and opportunity for a subtask-driven approach to RL through
44 the lens of CVaR optimization. In Section 5, we introduce the RED RL framework, including the
45 concept of the reward-extended TD error. We also introduce a family of RED RL algorithms for
46 prediction and control, and highlight their convergence properties (with full convergence proofs in
47 Appendix B). In Section 6, we use the RED RL framework to derive a subtask-driven approach for
48 CVaR optimization, and provide empirical results which show that this approach can be used to suc-
49 cessfully learn a policy that optimizes the CVaR risk measure. Finally, in Section 7, we emphasize
50 our framework's potential usefulness towards tackling other challenging problems outside the realm
51 of risk-awareness, highlight some of its limitations, and suggest some directions for future research.

## 2 Related Work

53 **Average-Reward Reinforcement Learning:** Average-reward (or average-cost) MDPs, despite be-
54 ing one of the most well-studied frameworks for sequential decision-making under uncertainty (Put-
55 erman, 1994), have remained relatively unexplored in reinforcement learning (RL) settings. To date,
56 notable works on the subject (in the context of RL) include Schwartz (1993), Tsitsiklis and Van Roy
57 (1999), Abounadi et al. (2001), Gosavi (2004), Bhatnagar et al. (2009), and Wan et al. (2021). Most
58 relevant to our work is Wan et al. (2021), which provided a rigorous theoretical treatment of average-
59 reward MDPs in the context of RL, and proposed the proven-convergent 'Differential Q-learning'
60 and 'Differential TD-learning' algorithms. Our work builds on the methods from Wan et al. (2021)
61 to develop a theoretical framework for solving various learning objectives simultaneously.

62 We note that these learning objectives, or *subtasks*, as explored in our work, are different to that of
63 hierarchical RL (e.g. Sutton et al. (1999)). In particular, in hierarchical RL, the focus is on using
64 temporally-abstracted actions, known as 'options' (or 'skills'), such that the agent learns a policy for
65 each option, as well as an inter-option policy. By contrast, in our work we learn a single policy, and
66 the subtasks are not part of the action-space. Similarly, the notion of solving multiple objectives in
67 parallel has been widely-explored in the discounted setting (e.g. McLeod et al. (2021)). However,
68 much of this work focuses on learning multiple state representations (or 'features'), options, policies,
69 and/or value functions. By contrast, in our work we learn a single policy and value function, and the
70 subtasks are not part of the state or action-spaces. To the best of our knowledge, our work is the first
71 to explore solving subtasks simultaneously in the average-reward setting.

72 **Risk-Aware Learning and Optimization in MDPs:** The notion of risk-aware learning and opti-
73 mization in MDP-based settings has been long-studied, from the well-established expected utility
74 framework (Howard and Matheson, 1972), to the more contemporary framework of coherent risk
75 measures (Artzner et al., 1999). To date, these risk-based efforts have almost exclusively focused
76 on the discounted setting. Importantly, optimizing the CVaR risk measure in these settings typi-
77 cally requires augmenting the state-space and/or having to utilize an explicit bi-level optimization
78 scheme, which can, for example, involve solving multiple MDPs. Seminal works that have looked at
79 CVaR optimization in the standard discounted setting include Bäuerle and Ott (2011) and Chow et al.
80 (2015); Hau et al. (2023a). In the distributional setting, works such as Dabney et al. (2018) have pro-
81 posed a CVaR optimization approach that does not require an augmented state-space or an explicit
82 bi-level optimization, however it was later shown by Lim and Malik (2022) that such an approach
83 converges to neither the optimal dynamic-CVaR nor the optimal static-CVaR policies (Lim and Ma-
84 lik (2022) then proposed a valid approach that utilizes an augmented state-space). Some works have
85 looked at optimizing a time-consistent (Ruszczyński, 2010) interpretation of CVaR, however this

86  only approximates CVaR, as CVaR is not a time-consistent risk measure (Boda and Filar, 2006).
87  Other works have looked at optimizing similar objectives to CVaR that are more computationally
88  tractable, such as the entropic value-at-risk (Hau et al., 2023b).

89  Most similar to our work (in non-average-reward settings) is Stanko and Macek (2019), where the
90  authors used a vaguely similar update to the one derived in our work. However, all of the methods
91  proposed in Stanko and Macek (2019) require either an augmented state-space or an explicit bi-
92  level optimization. In the average-reward setting, Xia et al. (2023) proposed a set of algorithms
93  for optimizing the CVaR risk measure, however their methods require the use of an augmented
94  state-space and a sensitivity-based bi-level optimization. By contrast, our work, to the best of our
95  knowledge, is the first to optimize CVaR in an MDP-based setting without the use of an explicit
96  bi-level optimization scheme or an augmented state-space. We note that other works have looked at
97  optimizing other risk measures in the average-reward setting, such as the exponential cost (Murthy
98  et al., 2023), and variance (Prashanth and Ghavamzadeh, 2016).

## 3 Preliminaries

### 3.1 Average-Reward Reinforcement Learning

101  A finite average-reward MDP is the tuple $\mathcal{M} \doteq \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, p \rangle$, where $\mathcal{S}$ is a finite set of states, $\mathcal{A}$ is
102  a finite set of actions, $\mathcal{R} \subset \mathbb{R}$ is a bounded set of rewards, and $p : \mathcal{S} \times \mathcal{A} \times \mathcal{R} \times \mathcal{S} \to [0, 1]$ is
103  a probabilistic transition function that describes the dynamics of the environment. At each discrete
104  time step, $t = 0, 1, 2, \ldots$, an agent chooses an action, $A_t \in \mathcal{A}$, based on its current state, $S_t \in \mathcal{S}$,
105  and receives a reward, $R_{t+1} \in \mathcal{R}$, while transitioning to a (potentially) new state, $S_{t+1}$, such that
106  $p(s', r \mid s, a) = \mathbb{P}(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$. In an average-reward MDP, an agent
107  aims to find a policy, $\pi : \mathcal{S} \to \mathcal{A}$, that optimizes the long-run (or limiting) average-reward, $\bar{r}$, which
108  is defined as follows for a given policy, $\pi$:

$$\bar{r}_\pi(s) \doteq \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t \mid S_0 = s, A_{0:t-1} \sim \pi]. \tag{1}$$

109  In this work, we limit our discussion to *stationary Markov* policies, which are time-independent
110  policies that satisfy the Markov property.

111  When working with average-reward MDPs, it is common to simplify Equation (1) into a more work-
112  able form by making certain assumptions about the Markov chain induced by following policy $\pi$. To
113  this end, a *unichain* assumption is typically used when doing prediction (learning) because it ensures
114  the existence of a unique limiting distribution of states, $\mu_\pi(s) \doteq \lim_{t \to \infty} \mathbb{P}(S_t = s \mid A_{0:t-1} \sim \pi)$,
115  that is independent of the initial state, thereby simplifying Equation (1) to the following:

$$\bar{r}_\pi = \sum_{s \in \mathcal{S}} \mu_\pi(s) \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r \mid s, a) r. \tag{2}$$

116  Similarly, a *communicating* assumption is typically used for control (optimization) because it en-
117  sures the existence of a unique optimal average-reward, $\bar{r}*$, that is independent of the initial state.

118  To solve an average-reward MDP, solution methods such as dynamic programming or RL can be
119  used in conjunction with the following *Bellman* (or *Poisson*) equations:

$$v_\pi(s) = \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a)[r - \bar{r}_\pi + v_\pi(s')], \tag{3}$$

$$q_\pi(s, a) = \sum_{s'} \sum_r p(s', r \mid s, a)[r - \bar{r}_\pi + \max_{a'} q_\pi(s', a')], \tag{4}$$

120  where, $v_\pi(s)$ is the state-value function and $q_\pi(s, a)$ is the state-action value function for a given
121  policy, $\pi$. Solution methods for average-reward MDPs are typically referred to as *differential* meth-
122  ods because of the reward difference (i.e., $r - \bar{r}_\pi$) operation that occurs in Equations (3) and (4). We

123 note that solution methods typically find solutions to Equations (3) and (4) up to a constant, $c$. This
124 is typically not a concern, given that the relative ordering of policies is usually what is of interest.

125 In the context of RL, Wan et al. (2021) proposed the tabular 'Differential TD-learning' and 'Dif-
126 ferential Q-learning' algorithms, which are able to learn and/or optimize the value function and
127 average-reward simultaneously using only the TD error. The 'Differential TD-learning' algorithm
128 is shown below:

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha_t \rho_t \delta_t \tag{5a}$$

$$V_{t+1}(s) \doteq V_t(s), \quad \forall s \neq S_t \tag{5b}$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t) \tag{5c}$$

$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \rho_t \delta_t \tag{5d}$$

129 where, $V_t : \mathcal{S} \to \mathbb{R}$ is a table of state-value function estimates, $\alpha_t$ is the step size, $\delta_t$ is the TD error,
130 $\rho_t \doteq \pi(A_t \mid S_t) / B(A_t \mid S_t)$ is the importance sampling ratio (with behavior policy, $B$), $\bar{R}_t$ is an
131 estimate of the average-reward, $\bar{r}_\pi$, and $\eta$ is a positive scalar.

## 3.2 Conditional Value-at-Risk (CVaR)

133 Consider a random variable $X$ with a finite mean on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, and with a
134 cumulative distribution function $F(x) = \mathbb{P}(X \leq x)$. The (left-tail) *value-at-risk (VaR)* of $X$ with
135 parameter $\tau \in (0, 1)$ represents the $\tau$-quantile of $X$, such that $\text{VaR}_\tau(X) = \sup\{x \mid F(x) \leq \tau\}$.
136 The (left-tail) *conditional value-at-risk (CVaR)* of $X$ with parameter $\tau$ is defined as follows:

$$\text{CVaR}_\tau(X) = \frac{1}{\tau} \int_0^\tau \text{VaR}_u(X) du. \tag{6}$$

137 When $F(X)$ is continuous at $x = \text{VaR}_\tau(X)$, $\text{CVaR}_\tau(X)$ can be interpreted as the expected value of
138 the $\tau$ left quantile of the distribution of $X$, such that $\text{CVaR}_\tau(X) = \mathbb{E}[X \mid X \leq \text{VaR}_\tau(X)]$.

139 Importantly, CVaR can be formulated as the following optimization (Rockafellar and Uryasev,
140 2000):

$$\text{CVaR}_\tau(X) = \sup_{y \in \mathbb{R}} \mathbb{E}[y - \frac{1}{\tau}(y - X)^+] = \mathbb{E}[\text{VaR}_\tau(X) - \frac{1}{\tau}(\text{VaR}_\tau(X) - X)^+], \tag{7}$$

141 where, $(u)^+ = \max(u, 0)$. Existing MDP-based methods typically leverage the above formulation
142 when optimizing for CVaR, by augmenting the state-space with a state that corresponds (either
143 directly or indirectly) to an estimate of $\text{VaR}_\tau(X)$ (in this case, $y$), and solving the following bi-level
144 optimization:

$$\sup_\pi \text{CVaR}_\tau(X) = \sup_\pi \sup_{y \in \mathbb{R}} \mathbb{E}[y - \frac{1}{\tau}(y - X)^+] = \sup_{y \in \mathbb{R}}(y - \frac{1}{\tau} \sup_\pi \mathbb{E}[(y - X)^+]), \tag{8}$$

145 where the 'inner' optimization problem can be solved using standard MDP solution methods.

146 In discounted MDPs, the random variable $X$ corresponds to a (potentially-discounted) sum of re-
147 wards. In average-reward MDPs, $X$ corresponds to the limiting per-step reward. In other words,
148 the natural interpretation of CVaR in the average-reward setting is that of the CVaR of the limiting
149 reward distribution, as shown below (for a given policy, $\pi$) (Xia et al., 2023):

$$\text{CVaR}_{\tau,\pi}(s) \doteq \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^n \text{CVaR}_\tau[R_t \mid S_0 = s, A_{0:t-1} \sim \pi]. \tag{9}$$

150 As with the average-reward (i.e., Equation (1)), a unichain assumption (or similar) makes this CVaR
151 objective independent of the initial state. In recent years, CVaR has emerged as a popular risk

152  measure, in-part because it is a 'coherent' risk measure (Artzner et al., 1999), meaning that it satisfies
153  key mathematical properties which can be meaningful in safety-critical and risk-related applications.

154  Figure 1 depicts the agent-environment interaction in an average-reward MDP, where following
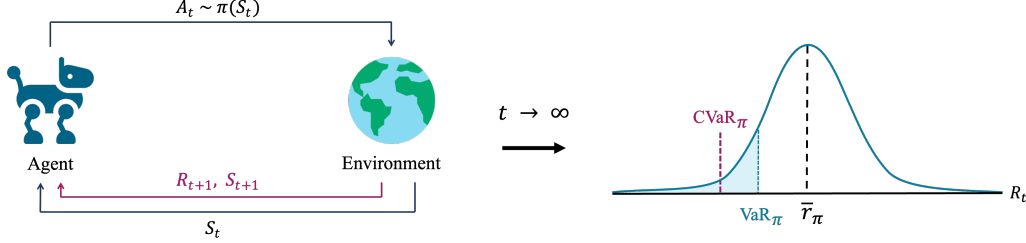155  policy $\pi$ yields a limiting average-reward and reward CVaR.



Figure 1: Illustration of the agent-environment interaction in an average-reward MDP. As $t \to \infty$, following policy $\pi$ yields a limiting per-step reward distribution with an average-reward, $\bar{r}_\pi$, and a conditional value-at-risk, $\text{CVaR}_\pi$. Standard average-reward RL methods aim to optimize the average-reward, $\bar{r}_\pi$. By contrast, in our work we aim to optimize $\text{CVaR}_\pi$.

## 4   A Subtask-Driven Approach

157  In this section, we motivate the need and opportunity for a subtask-driven approach to RL through
158  the lens of CVaR optimization. Let us begin by considering the standard approach used by existing
159  MDP-based methods for CVaR optimization. This approach, which is described in Equation (8),
160  requires that we pick a wide range of guesses for the optimal value-at-risk, VaR, and that for each
161  guess, $y$, we solve an MDP. Then, out of all of the MDP solutions, we pick the best one as our final
162  solution (which corresponds to $y = \text{VaR}$). Moreover, to further compound the computational costs,
163  this approach requires that the state-space be augmented with a state that corresponds (either directly
164  or indirectly) to the VaR guess, $y$ (e.g. see Bäuerle and Ott (2011)). Hence, this approach requires
165  the use of both an explicit bi-level optimization scheme, and an augmented state-space. Importantly
166  however, this computationally-expensive process would not be needed if we somehow knew what
167  the optimal value for $y$ (i.e., VaR) was. In fact, in the average-reward setting, if we know this optimal
168  value, VaR, then optimizing for CVaR ultimately amounts to optimizing an average (as per Equation
169  (7)), which can be done trivially using the standard average-reward MDP.

170  As such, it would appear that, to optimize CVaR, we are stuck between two extremes: a significantly
171  computationally-expensive process if we do not know the optimal value-at-risk, VaR, and a trivial
172  process if we do. But what if we could estimate VaR along the way? That is, keep some sort of
173  running estimate of VaR that we optimize simultaneously as we optimize CVaR. Indeed, such an
174  approach has been proposed in the discounted setting (e.g. Stanko and Macek (2019)), however,
175  no approach has been able to successfully remove both the augmented state-space and the explicit
176  bi-level optimization requirements. The primary difficulty lies in *how* one updates the estimate of
177  VaR along the way.

178  Critically, this is where the findings from Wan et al. (2021) come into play. In particular, Wan et al.
179  (2021) proposed proven-convergent algorithms for the average-reward setting that can learn and/or
180  optimize the value function and average-reward simultaneously using only the TD error. In other
181  words, these algorithms are able to solve two learning objectives simultaneously using only the TD
182  error. Yet, the focus in Wan et al. (2021) was on proving the convergence of such algorithms, without
183  exploring the underlying structural properties of the average-reward MDP that made such a process
184  possible to begin with. In this work, we formalize these underlying properties, and utilize them to
185  show that if one modifies, or *extends*, the reward from the MDP with various learning objectives
186  that satisfy certain key properties, then these objectives, or *subtasks*, can be solved simultaneously
187  using a modified, or *reward-extended*, version of the TD error. Consequently, in terms of CVaR

188  optimization, this allows us to develop appropriate learning updates for the VaR and CVaR estimates
189  based solely on the TD error, such that we no longer need to augment the state-space or perform an
190  explicit bi-level optimization.

191  In Section 5, we present the theoretical framework that enables the aforementioned subtask-driven
192  approach. Then, in Section 6, we adapt this general-purpose framework for CVaR optimization.

## 5  Reward-Extended Differential (RED) Reinforcement Learning

194  In this section, we present our primary contribution: a framework for solving various learning ob-
195  jectives, or *subtasks*, simultaneously in the average-reward setting. We call this framework *reward-*
196  *extended differential* (or *RED*) reinforcement learning. The 'differential' part of the name comes
197  from the use of the differential algorithms from average-reward MDPs. The 'reward-extended' part
198  of the name comes from the use of the *reward-extended TD error*, a novel concept that we will
199  introduce shortly. Through this framework, we show how the average-reward MDP's unique struc-
200  tural properties can be leveraged to solve (i.e., learn or optimize) any given subtask using only a
201  TD error-based update. We first provide a formal definition for a (generic) subtask, then proceed
202  to derive a framework that allows us to solve any given subtask that satisfies this definition. In the
203  subsequent section, we utilize this framework to tackle the CVaR optimization problem.

204  **Definition 5.1** (Subtask). *A subtask, $z_i$, is any scalar prediction or control objective belonging to*
205  *a corresponding bounded set $\mathcal{Z}_i \subset \mathbb{R}$, such that there exists a linear or piecewise linear subtask*
206  *function, $f : \mathcal{R} \times \mathcal{Z}_1 \times \mathcal{Z}_2 \times \cdots \times \mathcal{Z}_i \times \cdots \times \mathcal{Z}_n \to \tilde{\mathcal{R}}$, where $\mathcal{R}$ is the bounded set of observed*
207  *per-step rewards from the MDP $\mathcal{M}$, $\tilde{\mathcal{R}} \subset \mathbb{R}$ is a bounded set of 'extended' per-step rewards whose*
208  *long-run average is the primary prediction or control objective of the MDP, $\tilde{\mathcal{M}} \doteq \langle \mathcal{S}, \mathcal{A}, \tilde{\mathcal{R}}, \tilde{p} \rangle$, and*
209  *$\mathcal{Z} = \{z_1 \in \mathcal{Z}_1, z_2 \in \mathcal{Z}_2, \ldots, z_n \in \mathcal{Z}_n\}$ is the set of $n$ subtasks that we wish to solve, such that:*

210  *i) $f$ is invertible with respect to each input given all other inputs; and*

211  *ii) each subtask $z_i \in \mathcal{Z}$ in $f$ is independent of the states and actions, and hence independent of*
212  *the observed per-step reward, $R_t \in \mathcal{R}$, such that $\mathbb{P}(S_{t+1} = s', \tilde{R}_{t+1} = f(r, z_1, \ldots, z_n) \mid S_t =$*
213  *$s, A_t = a) = \mathbb{P}(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$, and $\mathbb{E}[f_j(R_t, z_1, z_2, \ldots, z_n)] =$*
214  *$f_j(\mathbb{E}[R_t], z_1, z_2, \ldots, z_n)$, where $f_j$ denotes the jth segment of a piecewise linear subtask function,*
215  *and $\mathbb{E}$ denotes any expectation taken with respect to the states and actions.*

216  In essence, the above definition states that a subtask is some constant, $z_i$, that we wish to learn and/or
217  optimize. From an algorithmic perspective, this means that we will start with some initial estimate
218  (or guess) for the subtask, $Z_{i,t}$, then update this estimate at every time step, such that $Z_{i,t} \to z_i$ or
219  $Z_{i,t} \to z_i^*$, depending on whether we are doing prediction or control (where $z_i^*$ denotes the optimal
220  subtask value). But how can we derive an appropriate update rule that accomplishes this? In the
221  following section, we will introduce the *reward-extended TD error*, through which we can derive
222  such an update rule for any subtask that satisfies Definition 5.1, such that $Z_{i,t} \to z_i$ when doing
223  prediction and $Z_{i,t} \to z_i^*$ when doing control.

### 5.1  The Reward-Extended TD Error

225  In this section, we introduce and derive the *reward-extended TD error*. In particular, we derive
226  a generic, subtask-specific, TD-like error, $\beta_{i,t}$, through which we can learn and/or optimize any
227  subtask that satisfies Definition 5.1 via the update rule: $Z_{i,t+1} = Z_{i,t} + \eta\alpha_t\beta_{i,t}$, where $Z_{i,t}$ is an
228  estimate of subtask $z_i$, $\eta\alpha_t$ is the step size, and $\beta_{i,t}$ is the reward-extended TD error for subtask $z_i$.

229  Importantly, we will show that the reward-extended TD error satisfies the following property:
230  $\mathbb{E}_\pi[\beta_{i,t}] \to 0 \ \forall i = 1, 2, \ldots, n$ as $\mathbb{E}_\pi[\delta_t] \to 0$, where $\delta_t$ is the regular TD error, such that min-
231  imizing the regular TD error allows us to solve all subtasks simultaneously. This motivates our
232  naming of the reward-extended TD error, given that it is intrinsically tied to the regular TD error.

6

Let us begin by considering the common RL update rule of the form: *NewEstimate ← OldEstimate + StepSize [Target − OldEstimate]* (Sutton and Barto, 2018; Naik, 2024). Our aim is to find an appropriate set of subtask-specific 'targets', $\{\phi_{i,t}\}_{i=1}^n$, such that $\mathbb{E}_\pi[\beta_{i,t}] = \mathbb{E}_\pi[\phi_{i,t} - Z_{i,t}] \to 0 \ \forall i = 1, 2, \ldots, n$ as $\mathbb{E}_\pi[\delta_t] \to 0$. To this end, let us consider a generic piecewise linear subtask function with $m$ piecewise segments:

$$
\tilde{R}_t = \begin{cases} b_r^1 R_t + b_0^1 + b_1^1 z_1 + b_2^1 z_2 + \ldots + b_n^1 z_n, & r_0 \le R_t < r_1 \\ b_r^2 R_t + b_0^2 + b_1^2 z_1 + b_2^2 z_2 + \ldots + b_n^2 z_n, & r_1 \le R_t < r_2 \\ \vdots \\ b_r^m R_t + b_0^m + b_1^m z_1 + b_2^m z_2 + \ldots + b_n^m z_n, & r_{m-1} \le R_t \le r_m \end{cases}, \tag{10}
$$

where $r_k \in \mathcal{R} \ \forall k = 0, 1, \ldots, m$, such that $r_0, r_m$ represent the lower and upper bounds of the observed per-step reward, $R_t$, respectively, $b_r^j, b_0^j \in \mathbb{R}$, and $b_i^j \in \mathbb{R} \setminus \{0\}$, where $b^j$ denotes a (predefined) constant in the $j$th segment of the piecewise linear subtask function.

Now, let us consider the TD error, $\delta_t$, associated with (10) in the prediction setting. Let $\tilde{R}_{j,t}$ be shorthand for the $j$th segment of (10), such that the TD error at any time step can be expressed as:

$$
\delta_{j,t} = \tilde{R}_{j,t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t) \tag{11a}
$$

$$
= b_r^j R_{t+1} + b_0^j + b_1^j Z_{1,t} + b_2^j Z_{2,t} + \ldots + b_n^j Z_{n,t} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t), \tag{11b}
$$

where $V_t : \mathcal{S} \to \mathbb{R}$ denotes a table of state-value function estimates, $\bar{R}_t$ denotes an estimate of the average-reward, $\bar{r}_\pi$, $Z_{i,t}$ denotes an estimate of subtask $z_i \ \forall i = 1, 2, \ldots, n$, and $j$ corresponds to the piecewise condition, $r_{j-1} \le R_{t+1} \le r_j$, that is satisfied by the observed per-step reward, $R_{t+1}$.

Hence, as learning progresses, different $\tilde{R}_{j,t+1}$ values will be used to define the TD error based on which piecewise condition is satisfied at a given time step. Moreover, we know that the probability that $\delta_t = \delta_{j,t}$ is equal to the probability that $r_{j-1} \le R_{t+1} < r_j$. This allows us to express the expected TD error associated with (10) as follows:

$$
\mathbb{E}_\pi[\delta_t] = \sum_{j=1}^m \mathbb{P}(r_{j-1} \le R_{t+1} < r_j)\mathbb{E}_\pi[\delta_{j,t}]. \tag{12}
$$

Now, let us consider the implications of $\mathbb{E}_\pi[\delta_t] \to 0$ as it relates to $\mathbb{E}_\pi[\delta_{j,t}]$. One possibility is that $\mathbb{E}_\pi[\delta_{j,t}] \to 0 \ \forall j = 1, 2, \ldots, m$. However, this may not necessarily be the case; it is possible that, for example, a pair of non-zero $\mathbb{P}(r_{j-1} \le R_{t+1} < r_j)\mathbb{E}_\pi[\delta_{j,t}]$ terms cancel each other out, such that $\mathbb{E}_\pi[\delta_t] \to 0$ but $\mathbb{E}_\pi[\delta_{j,t}] \to \lambda_j \ \forall j = 1, 2, \ldots, m$, where $\lambda_j \in \mathbb{R}$. In such a case, what we do know is that if $\mathbb{E}_\pi[\delta_t] \to 0$, then the Bellman equation (3) must be satisfied, such that: $V_t(s) = \mathbb{E}_\pi[\tilde{R}_{t+1} - \bar{R}_t + V_t(S_{t+1}) \mid S_t = s]$. As such, we can write the following expression for $\lambda_j$, and solve for an arbitrary subtask, $z_i$, as follows:

$$
\lambda_j = \mathbb{E}_\pi[\tilde{R}_{j,t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t)] \tag{13a}
$$

$$
= \mathbb{E}_\pi\left[\tilde{R}_{j,t+1} - \bar{R}_t + V_t(S_{t+1}) - \left(\tilde{R}_{t+1} - \bar{R}_t + V_t(S_{t+1})\right)\right] \tag{13b}
$$

$$
= \mathbb{E}_\pi[\tilde{R}_{j,t+1}] - \mathbb{E}_\pi[\tilde{R}_{t+1}] \tag{13c}
$$

$$
= \mathbb{E}_\pi[\tilde{R}_{j,t+1}] - \bar{r}_\pi \quad \text{(See Remark 5.3)} \tag{13d}
$$

$$
= \mathbb{E}_\pi[b_r^j R_{t+1} + b_0^j + \ldots + b_{i-1}^j z_{i-1} + b_{i+1}^j z_{i+1} + \ldots + b_n^j z_n - \bar{r}_\pi] + b_i^j z_i \tag{13e}
$$

$$
\implies z_i = \mathbb{E}_\pi\left[-\frac{1}{b_i^j}\left(b_r^j R_{t+1} + b_0^j + \ldots + b_{i-1}^j z_{i-1} + b_{i+1}^j z_{i+1} + \ldots + b_n^j z_n - \bar{r}_\pi - \lambda_j\right)\right] \tag{13f}
$$

$$
\doteq \mathbb{E}_\pi[\phi_{i,j}], \tag{13g}
$$

where we used the fact that $z_i$ is independent of the states and actions to pull it out of the expectation. Here, we use $\phi_{i,j}$ to denote the expression inside the expectation in Equation (13f).

Hence, to learn $z_i$ from experience, we can utilize the common RL update rule, using the term inside the expectation in Equation (13g), $\phi_{i,j}$, as the 'target', which yields the update:

$$Z_{i,t+1} = Z_{i,t} + \eta\alpha_t \begin{cases} \phi_{i,1,t} - Z_{i,t}, & r_0 \le R_{t+1} < r_1 \\ \vdots \\ \phi_{i,m,t} - Z_{i,t}, & r_{m-1} \le R_{t+1} \le r_m \end{cases} \tag{14a}$$

$$= Z_{i,t} + \eta\alpha_t \begin{cases} (-1/b_i^1)\left(\tilde{R}_{1,t+1} - \bar{R}_t - \delta_t\right), & r_0 \le R_{t+1} < r_1 \\ \vdots \\ (-1/b_i^m)\left(\tilde{R}_{m,t+1} - \bar{R}_t - \delta_t\right), & r_{m-1} \le R_{t+1} \le r_m \end{cases} \tag{14b}$$

$$\doteq Z_{i,t} + \eta\alpha_t\beta_{i,t}, \tag{14c}$$

where $Z_{i,t}$ is the estimate of subtask $z_i$ at time $t$, $\phi_{i,j,t} \doteq (-1/b_i^j)(b_r^j R_{t+1} + b_0^j + \ldots + b_{i-1}^j Z_{i-1,t} + b_{i+1}^j Z_{i+1,t} + \ldots + b_n^j Z_{n,t} - \bar{R}_t - \delta_t)$, and $\eta\alpha_t$ is the step size.

As such, we now have an expression for the reward-extended TD error for subtask $z_i$, $\beta_{i,t}$. We will now show that this term satisfies the desired property: $\mathbb{E}_\pi[\beta_{i,t}] \to 0 \; \forall i = 1, 2, \ldots, n$ as $\mathbb{E}_\pi[\delta_t] \to 0$, such that minimizing the regular TD error allows us to solve all the subtasks simultaneously:

**Theorem 5.1.** *Consider an average-reward MDP with a set of reward-extended TD errors, $\{\beta_{i,t}\}_{i=1}^n$, as defined in Equation (14), corresponding to a subtask function with $n$ subtasks that satisfy Definition 5.1. The set of reward-extended TD errors, $\{\beta_{i,t}\}_{i=1}^n$, satisfies the following property: $\mathbb{E}_\pi[\beta_{i,t}] \to 0 \; \forall i = 1, 2, \ldots, n$ as $\mathbb{E}_\pi[\delta_t] \to 0$, where $\beta_{i,t}$ denotes the reward-extended TD error for subtask $z_i$, and $\delta_t$ denotes the regular TD error.*

*Proof.* Let us consider the reward-extended TD error associated with an arbitrary $j$th segment of the piecewise linear subtask function for an arbitrary $i$th subtask: $\beta_{i,j,t} \doteq (-1/b_i^j)(\tilde{R}_{j,t+1} - \bar{R}_t - \delta_t)$. As $\mathbb{E}_\pi[\delta_t] \to 0$, $\bar{R}_t \to \bar{r}_\pi$ (by Theorem 3 of Wan et al. (2021)) and $\delta_t \to \lambda_j$ for this $j$th segment. Hence, $\mathbb{E}_\pi[\beta_{i,j,t}] \to (-1/b_i^j)(\mathbb{E}_\pi[\tilde{R}_{j,t+1}] - \bar{r}_\pi - \lambda_j) = (-1/b_i^j)(\lambda_j - \lambda_j) = 0$. Now, because we chose $j$ arbitrarily, we have, for all $j \in \{1, 2, \ldots, m\}$, that $\mathbb{E}_\pi[\beta_{i,j,t}] \to 0$. As such, and because we chose $i$ arbitrarily, we can conclude that $\mathbb{E}_\pi[\beta_{i,t}] = \sum_{j=1}^m \mathbb{P}(r_{j-1} \le R_{t+1} < r_j)\mathbb{E}_\pi[\beta_{i,j,t}] \to 0 \; \forall i = 1, 2, \ldots, n$ as $\mathbb{E}_\pi[\delta_t] \to 0$. This completes the proof. $\qquad\square$

As such, we have derived the desired update rule that we can use to solve any given subtask in the prediction setting. The same logic can be applied in the control setting to derive equivalent updates, where we note that it directly follows from Definition 5.1 that the existence of an optimal average-reward, $\bar{r}*$, implies the existence of corresponding optimal subtask values, $z_i^* \; \forall z_i \in \mathcal{Z}$.

**Remark 5.1.** *In the case of a (non-piecewise) linear subtask function, the expression for the reward-extended TD error can be simplified to $\beta_{i,t} \doteq (-1/b_i)\delta_t$ by setting $\lambda = 0$ in Equation (13a), solving for the target, $z_i$, and applying a similar process to the one described in Equation (14).*

**Remark 5.2.** *Given Remark 5.1, it can be shown that if one treats the average-reward, $\bar{r}_\pi$, as a subtask, and derives the reward-extended TD error for it, the process yields the average-reward update (e.g. Equation (5d)) from the Differential algorithms proposed in Wan et al. (2021). Hence, our work can be viewed as a generalization of the work performed in Wan et al. (2021).*

**Remark 5.3.** *Strictly speaking, $\bar{r}_\pi = \mathbb{E}_\pi[\tilde{R}_{t+1}] + c, \; c \in \mathbb{R}$. This is because average-reward solution methods typically find the solutions to the Bellman equations (3) and (4) up to an additive constant, c. This means that, like the average-reward estimate, our subtask estimates converge to the actual subtask values, up to an additive constant. For simplicity, we omit this additive constant in our work, unless strictly necessary, given that it is commonplace to assume that solutions in the average-reward setting are correct up to an additive constant.*

### 5.2 The RED Algorithms

In this section, we introduce the *RED RL algorithms*, which integrate the update rules derived in the previous section into the average-reward RL framework from Wan et al. (2021). The full algorithms, including algorithms that utilize function approximation, are included in Appendix A.

**RED TD-learning algorithm (tabular):** We update a table of estimates, $V_t : \mathcal{S} \to \mathbb{R}$ as follows:

$$\tilde{R}_{t+1} = f(R_{t+1}, Z_{1,t}, Z_{2,t}, \dots, Z_{n,t}) \tag{15a}$$

$$\delta_t = \tilde{R}_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t) \tag{15b}$$

$$V_{t+1}(S_t) = V_t(S_t) + \alpha_t \rho_t \delta_t \tag{15c}$$

$$\bar{R}_{t+1} = \bar{R}_t + \eta_r \alpha_t \rho_t \delta_t \tag{15d}$$

$$Z_{i,t+1} = Z_{i,t} + \eta_{z_i} \alpha_t \rho_t \beta_{i,t}, \quad \forall z_i \in \mathcal{Z} \tag{15e}$$

where, $R_t$ is the observed reward, $Z_{i,t}$ is an estimate of subtask $z_i$, $\beta_{i,t}$ is the reward-extended TD error for subtask $z_i$, $\alpha_t$ is the step size, $\delta_t$ is the TD error, $\rho_t$ is the importance sampling ratio, $\bar{R}_t$ is an estimate of the long-run average-reward of $\tilde{R}_t$, $\bar{r}_\pi$, and $\eta_r, \eta_{z_i}$ are positive scalars.

Wan et al. (2021) showed for their Differential TD-learning algorithm that $R_t$ converges to $\bar{r}_\pi$, and $V_t$ converges to a solution of $v$ in Equation (3) for a given policy, $\pi$. We now provide an equivalent theorem for our RED TD-learning algorithm, which also shows that $Z_{i,t}$ converges to $z_{i,\pi} \; \forall z_i \in \mathcal{Z}$, where $z_{i,\pi}$ denotes the subtask value induced when following policy $\pi$:

**Theorem 5.2** (informal). *The RED TD-learning algorithm* (15) *converges, almost surely, $\bar{R}_t$ to $\bar{r}_\pi$, $Z_{i,t}$ to $z_{i,\pi} \; \forall z_i \in \mathcal{Z}$, and $V_t$ to a solution of $v$ in the Bellman Equation* (3), *up to an additive constant, $c$, if the following assumptions hold: 1) the Markov chain induced by the target policy, $\pi$, is unichain, 2) every state–action pair for which $\pi(a|s) > 0$ occurs an infinite number of times under the behavior policy, 3) the step sizes are decreased appropriately, 4) the ratio of the update frequency of the most-updated state to the least-updated state is finite, 5) the subtasks are in accordance with Definition 5.1, and 6) the subtask step sizes are decreased appropriately.*

**RED Q-learning algorithm (tabular):** We update $Q_t : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ as follows:

$$\tilde{R}_{t+1} = f(R_{t+1}, Z_{1,t}, Z_{2,t}, \dots, Z_{n,t}) \tag{16a}$$

$$\delta_t = \tilde{R}_{t+1} - \bar{R}_t + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t) \tag{16b}$$

$$Q_{t+1}(S_t, A_t) = Q_t(S_t, A_t) + \alpha_t \delta_t \tag{16c}$$

$$\bar{R}_{t+1} = \bar{R}_t + \eta_r \alpha_t \delta_t \tag{16d}$$

$$Z_{i,t+1} = Z_{i,t} + \eta_{z_i} \alpha_t \beta_{i,t}, \quad \forall z_i \in \mathcal{Z} \tag{16e}$$

where, $R_t$ is the observed reward, $Z_{i,t}$ is an estimate of subtask $z_i$, $\beta_{i,t}$ is the reward-extended TD error for subtask $z_i$, $\alpha_t$ is the step size, $\delta_t$ is the TD error, $\bar{R}_t$ is an estimate of the long-run average-reward of $\tilde{R}_t$, $\bar{r}_\pi$, and $\eta_r, \eta_{z_i}$ are positive scalars. Wan et al. (2021) showed for their Differential Q-learning algorithm that $R_t$ converges to $\bar{r}*$, and $Q_t$ converges to a solution of $q$ in Equation (4). We now provide an equivalent theorem for our RED Q-learning algorithm, which also shows that $Z_{i,t}$ converges to the corresponding optimal subtask value $z_i^* \; \forall z_i \in \mathcal{Z}$:

**Theorem 5.3** (informal). *The RED Q-learning algorithm* (16) *converges, almost surely, $\bar{R}_t$ to $\bar{r}*$, $Z_{i,t}$ to $z_i^* \; \forall z_i \in \mathcal{Z}$, $\bar{r}_{\pi_t}$ to $\bar{r}*$, $z_{i,\pi_t}$ to $z_i^* \; \forall z_i \in \mathcal{Z}$, and $Q_t$ to a solution of $q$ in the Bellman Equation* (4), *up to an additive constant, $c$, where $\pi_t$ is any greedy policy with respect to $Q_t$, if the following assumptions hold: 1) the MDP is communicating, 2) the solution of $q$ in* (4) *is unique up to a constant, 3) the step sizes are decreased appropriately, 4) all the state–action pairs are updated an infinite number of times, 5) the ratio of the update frequency of the most-updated state–action pair to the least-updated state–action pair is finite, 6) the subtasks are in accordance with Definition 5.1, and 7) the subtask step sizes are decreased appropriately.*

See Appendix B for the formal version of these theorems, along with the full convergence proofs.

## 6 Case Study: RED RL for CVaR Optimization

In this section, we present a case-study which illustrates how the subtask-driven approach that was derived in Section 5 can be used to successfully tackle the CVaR optimization problem, *without* the use of an explicit bi-level optimization scheme (as in Equation (8)), or an augmented state-space.

First, in order to leverage the RED RL framework for CVaR optimization, we need to derive a valid subtask function for CVaR that satisfies the requirements of Definition 5.1. It turns out that we can use Equation (7) as a basis for the subtask function. The details of the adaptation of Equation (7) into a subtask function are presented in Appendix C. Critically, as discussed in Appendix C, optimizing the long-run average of the *extended* reward ($\tilde{R}_t$) from this subtask function corresponds to optimizing the long-run CVaR of the *observed* reward ($R_t$). Hence, we can utilize CVaR-specific versions of the RED algorithms presented in Equations (15) and (16) (or their non-tabular equivalents) to optimize VaR and CVaR, such that CVaR corresponds to the primary control objective (i.e., the $\bar{r}_\pi$ that we want to optimize), and VaR is the (single) subtask. We call the resulting algorithms, the *RED CVaR algorithms*. These algorithms, which are shown in full in Appendix C, update CVaR in an analogous way to the average-reward (i.e., CVaR corresponds to $\bar{R}_t$ in Equations (15) or (16)), and update VaR using a VaR-specific version of Equation (15e) or (16e) as follows:

$$
\text{VaR}_{t+1} = \begin{cases} \text{VaR}_t + \eta\alpha_t \left( \delta_t + \text{CVaR}_t - \text{VaR}_t \right), & R_{t+1} \geq \text{VaR}_t \\ \text{VaR}_t + \eta\alpha_t \left( \left( \frac{\tau}{\tau-1} \right) \delta_t + \text{CVaR}_t - \text{VaR}_t \right), & R_{t+1} < \text{VaR}_t \end{cases}, \tag{17}
$$

where, $\text{VaR}_t$ and $\text{CVaR}_t$ are estimates of VaR and CVaR, $\eta\alpha_t$ is the step size, $\tau$ is the CVaR parameter, $\delta_t$ is the TD error, and $R_t$ is the observed reward. As such, we are able to optimize VaR and CVaR without the use of an explicit bi-level optimization scheme or an augmented state-space.

We now present empirical results obtained when applying the RED CVaR algorithms on two RL tasks. The full set of experimental details and results can be found in Appendix D.

The first task corresponds to a two-state environment that we created for the purposes of testing our RED CVaR algorithms. It is called the *red-pill blue-pill* task (see Appendix E), where at every time step an agent can take either a 'red pill', which takes them to the 'red world' state, or a 'blue pill', which takes them to the 'blue world' state. Each state has its own characteristic per-step reward distribution, and in this case, for a sufficiently low CVaR parameter, $\tau$, the red world state has a reward distribution with a lower (worse) mean but higher (better) CVaR compared to the blue world state. As such, this task allows us to answer the following question: *can the RED CVaR algorithms successfully get the agent to learn a policy that prioritizes optimizing the reward CVaR over the average-reward?* In particular, we would expect that the RED CVaR algorithms learn a policy that prefers to stay in the red world, and that the (risk-neutral) Differential algorithms (from Wan et al. (2021)) learn a policy that prefers to stay in the blue world. This task is illustrated in Figure 2.
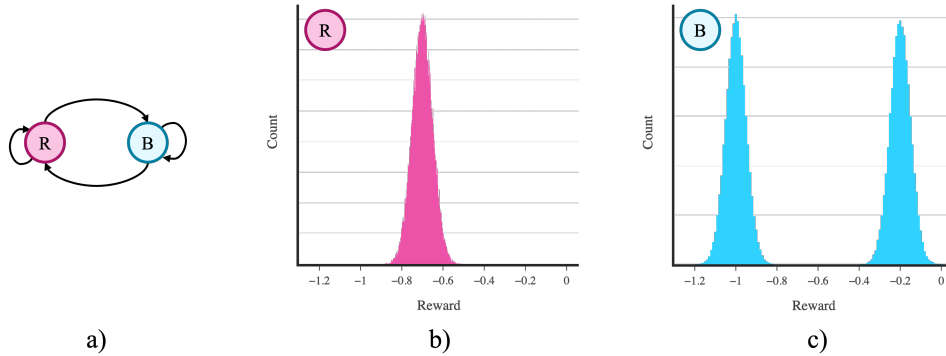


Figure 2: **a)** The *red-pill blue-pill* environment. **b) + c)** Histograms showing the per-step reward distribution of the **b)** 'red world', and **c)** 'blue world' states in the red-pill blue-pill environment.

365  The second task is the well-known *inverted pendulum* task, where an agent learns how to optimally
366  balance an inverted pendulum. We chose this task because it provides us with the opportunity to test
367  our algorithms in an environment where: 1) we must use function approximation (given the high-
368  dimensional state-space), and 2) where the optimal CVaR policy and the optimal average-reward
369  policy is the same policy (i.e., the policy that best balances the pendulum will yield a limiting re-
370  ward distribution with both the optimal average-reward and reward CVaR). This hence allows us to
371  directly compare the performance of our RED algorithms to that of the regular Differential algo-
372  rithms, as well as to gauge how function approximation affects the performance of our algorithms.
373  For this task, we utilized a simple actor-critic architecture (Barto et al., 1983; Sutton and Barto,
374  2018) as this allowed us to compare the performance of a (non-tabular) RED TD-learning algorithm
375  with a (non-tabular) Differential TD-learning algorithm.

376  In terms of empirical results, Figure 3 shows rolling averages of the average-reward and reward
377  CVaR as learning progresses in both tasks when using the regular Differential learning algorithms
378  (to optimize the average-reward) vs. the RED CVaR algorithms (to optimize the reward CVaR).
379  As shown in the figure, in the red-pill blue-pill task, the RED CVaR algorithm is able to success-
380  fully learn a policy that prioritizes maximizing the reward CVaR over the average-reward, thereby
381  achieving a sort of *risk-awareness*. In the inverted pendulum task, both methods converge to the
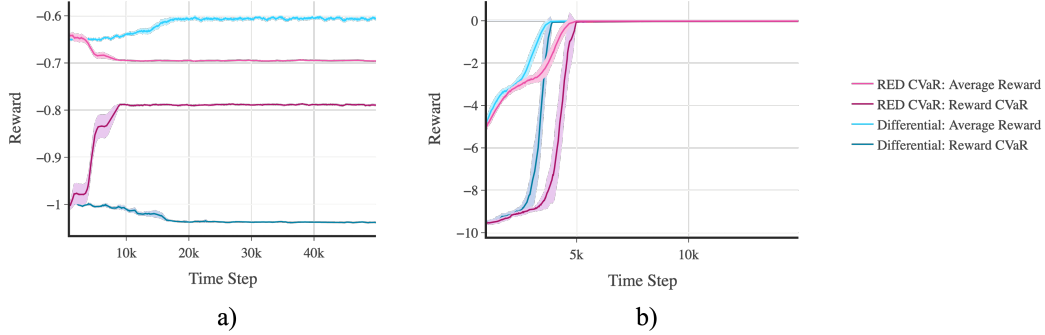382  same policy, as expected.



Figure 3: Rolling average-reward and reward CVaR as learning progresses when using the (risk-neutral) Differential algorithms vs. the (risk-aware) RED CVaR algorithms in the **a)** red-pill blue-pill, and **b)** inverted pendulum tasks. A solid line denotes the mean average-reward or reward CVaR, and the corresponding shaded region denotes a 95% confidence interval over **a)** 50 runs , or **b)** 10 runs. As shown in the figure, the RED CVaR algorithms are able to successfully learn a policy that prioritizes maximizing the reward CVaR, thereby achieving a sort of *risk-awareness*.
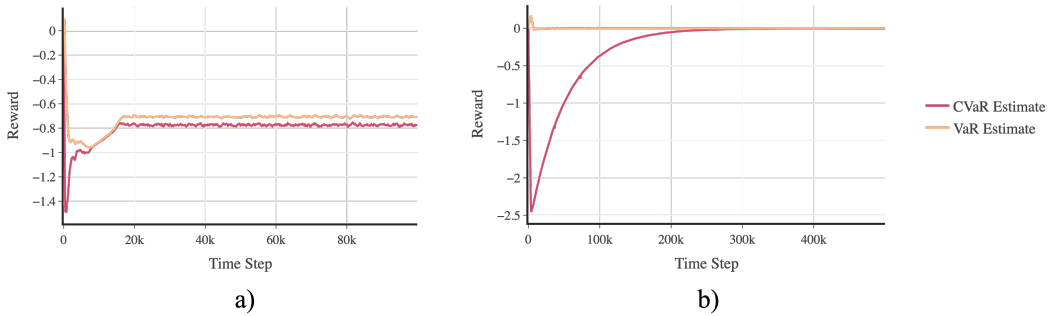


Figure 4: Typical convergence plots of the agent's VaR and CVaR estimates as learning progresses when using the RED CVaR algorithms in the **a)** red-pill blue-pill, and **b)** inverted pendulum tasks with an initial guess of 0.0 for both estimates.

11

Figure 4 shows typical convergence plots of the agent's VaR and CVaR estimates as learning progresses in both tasks when using the RED CVaR algorithms. As shown in the figure, the estimates converge in both tasks. In particular, the estimates converge to the correct VaR and CVaR values, up to an additive constant, thereby yielding the optimal CVaR policy, and hence, the results in Figure 3.

# 7 Discussion, Limitations, and Future Work

In this work, we introduced *reward-extended differential* (or *RED*) reinforcement learning: a novel reinforcement learning framework that can be used to solve various learning objectives, or *subtasks*, simultaneously in the average-reward setting. We introduced a family of RED RL algorithms for prediction and control, and then showcased how these algorithms could be utilized to effectively and efficiently tackle the CVaR optimization problem. More specifically, we were able to use the RED RL framework to derive a set of algorithms that can optimize the CVaR risk measure without using an explicit bi-level optimization scheme or an augmented state-space, thereby alleviating some of the computational challenges and non-trivialities that arise when performing risk-based optimization in the discounted setting. Empirically, we showed that the RED-based CVaR algorithms fared well in both tabular and linear function approximation settings.

More broadly, our work has introduced a theoretically-sound framework that allows for a subtask-driven approach to reinforcement learning, where various learning objectives (or subtasks) are solved simultaneously to help solve a larger, central learning objective. In this work, we showed (both theoretically and empirically) how this framework can be utilized to predict and/or optimize any arbitrary number of subtasks simultaneously in the average-reward setting. Central to this result is the novel concept of the reward-extended TD error, which is utilized in our framework to develop learning rules for the subtasks, and satisfies key theoretical properties that make it possible to solve any given subtask in a fully-online manner by minimizing the regular TD error. Moreover, we built upon existing results from Wan et al. (2021) to show the almost sure convergence of tabular algorithms derived from our framework. While we have only begun to grasp the implications of our framework, we have already seen some promising indications in the CVaR case study: the ability to turn explicit bi-level optimization problems into implicit bi-level optimizations that can be solved in a fully-online manner, as well as the potential to turn certain states (that meet certain conditions) into subtasks, thereby reducing the size of the state-space.

Nonetheless, while these results are encouraging, they are subject to a number of limitations. Firstly, by nature of operating in the average-reward setting, we are subject to the somewhat-strict assumptions made about the Markov chain induced by the policy (e.g. unichain or communicating). These assumptions could restrict the applicability of our framework, as they may not always hold in practice. Similarly, our definition for a subtask requires that the associated subtask function be linear or piecewise linear with respect to the subtasks, which may limit the applicability of our framework to simpler subtask functions. Finally, it remains to be seen empirically how our framework performs when dealing with multiple subtasks, when taking on more complex tasks, and/or when utilizing nonlinear function approximation.

Future work should look to address these limitations, as well as explore how these promising results can be extended to other domains, beyond the risk-awareness problem. In particular, we believe that the ability to optimize various subtasks simultaneously, as well as the potential to reduce the size of the state-space, by converting certain states to subtasks (where appropriate), could help alleviate significant computational challenges in other areas moving forward.

# References

J Abounadi, D Bertsekas, and V S Borkar. Learning algorithms for markov decision processes with average cost. *SIAM Journal on Control and Optimization*, 40(3):681–698, 2001.

Philippe Artzner, Freddy Delbaen, Jean-Marc Eber, and David Heath. Coherent measures of risk. *Math. Finance*, 9(3):203–228, July 1999.

Andrew G Barto, Richard S Sutton, and Charles W Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Trans. Syst. Man Cybern.*, SMC-13(5):834–846, September 1983.

Dimitri Bertsekas and John N Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.

Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor–critic algorithms. *Automatica (Oxf.)*, 45(11):2471–2482, November 2009.

Kang Boda and Jerzy A Filar. Time consistent dynamic risk measures. *Math. Methods Oper. Res. (Heidelb.)*, 63(1):169–186, February 2006.

Vivek S Borkar. Asynchronous stochastic approximations. *SIAM Journal on Control snd Optimization*, 36(3):840–851, 1998.

Vivek S Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Springer, 2009.

Nicole Bäuerle and Jonathan Ott. Markov decision processes with average-value-at-risk criteria. *Math. Methods Oper. Res.*, 74(3):361–379, December 2011.

Yinlam Chow, Aviv Tamar, Shie Mannor, and Marco Pavone. Risk-sensitive and robust decision-making: a CVaR optimization approach. In *Advances in Neural Information Processing Systems 28*, 2015.

Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for distributional reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.

Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Mach. Learn.*, 110(9):2419–2468, September 2021.

Abhijit Gosavi. Reinforcement learning for long-run average cost. *Eur. J. Oper. Res.*, 155(3):654–674, June 2004.

Jia Lin Hau, Erick Delage, Mohammad Ghavamzadeh, and Marek Petrik. On dynamic programming decompositions of static risk measures in markov decision processes. In *Advances in Neural Information Processing Systems 37*, 2023a.

Jia Lin Hau, Marek Petrik, and Mohammad Ghavamzadeh. Entropic risk optimization in discounted MDPs. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, 2023b.

Ronald A Howard and James E Matheson. Risk-sensitive markov decision processes. *Manage. Sci.*, 18(7):356–369, March 1972.

Shiau Hong Lim and Ilyas Malik. Distributional reinforcement learning for risk-sensitive policies. In *Advances in Neural Information Processing Systems 35*, 2022.

Matthew McLeod, Chunlok Lo, Matthew Schlegel, Andrew Jacobsen, Raksha Kumaraswamy, Marhta White, and Adam White. Continual auxiliary task learning. In *Advances in Neural Information Processing Systems 34*, 2021.

Yashaswini Murthy, Mehrdad Moharrami, and R Srikant. Modified policy iteration for exponential cost risk sensitive MDPs. In *Proceedings of the 5th Annual Learning for Dynamics and Control Conference*, 2023.

Abhishek Naik. *Reinforcement Learning for Continuing Problems Using Average Reward*. PhD thesis, University of Alberta, 2024.

L A Prashanth and Mohammad Ghavamzadeh. Variance-constrained actor-critic algorithms for discounted and average reward MDPs. *Mach. Learn.*, 105(3):367–417, December 2016.

Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 1994.

R Tyrrell Rockafellar and Stanislav Uryasev. Optimization of conditional value-at-risk. *The Journal of Risk*, 2(3):21–41, 2000.

Andrzej Ruszczyński. Risk-averse dynamic programming for markov decision processes. *Math. Program.*, 125(2):235–261, October 2010.

Anton Schwartz. A reinforcement learning method for maximizing undiscounted rewards. In *International Conference on Machine Learning*, 1993.

Silvestr Stanko and Karel Macek. Risk-averse distributional reinforcement learning: A CVaR optimization approach. In *Proceedings of the 11th International Joint Conference on Computational Intelligence*, 2019.

Richard S Sutton. Learning to predict by the methods of temporal differences. *Mach. Learn.*, 3:944, 1988.

Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction, 2nd edition*. MIT Press, November 2018.

Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artif. Intell.*, 112(1):181–211, August 1999.

John N Tsitsiklis and Benjamin Van Roy. Average cost temporal-difference learning. *Automatica*, 35:1799, 1999.

Yi Wan, Abhishek Naik, and Richard S Sutton. Learning and planning in average-reward markov decision processes. In *Proceedings of the 38th International Conference on Machine Learning*, 2021.

Li Xia, Luyao Zhang, and Peter W Glynn. Risk-sensitive markov decision processes with long-run CVaR criterion. *Prod. Oper. Manag.*, 32(12):4049–4067, December 2023.

## A  RED RL Algorithms

In this appendix, we provide pseudocode for our RED RL algorithms. We first present tabular algorithms, whose convergence proofs are included in Appendix B, and then provide equivalent algorithms that utilize function approximation.

---

**Algorithm 1** RED TD-Learning (Tabular)

---

**Input:** the policy $\pi$ to be evaluated, policy $B$ to be used, piecewise linear subtask function $f$ with $n$ subtasks, $m$ piecewise segments, piecewise conditions $r_{j-1} \leq R < r_j$ such that $f_j$ denotes the $j$th segment of $f$ that satisfies $r_{j-1} \leq R < r_j$, and constants $b_1^j, b_2^j, \ldots, b_n^j \ \forall j = 1, 2, \ldots, m$

**Algorithm parameters:** step size parameters $\alpha, \eta_r, \eta_{z_1}, \eta_{z_2}, \ldots, \eta_{z_n}$

Initialize $V(s) \ \forall s; \bar{R}$ arbitrarily (e.g. to zero)

Initialize subtasks $Z_1, Z_2, \ldots, Z_n$ arbitrarily (e.g. to zero)

Obtain initial $S$

**while** still time to train **do**

    $A \leftarrow$ action given by $B$ for $S$

    Take action $A$, observe $R, S'$

    $\tilde{R} = f(R, Z_1, Z_2, \ldots, Z_n)$

    $\delta = \tilde{R} - \bar{R} + V(S') - V(S)$

    $\rho = \pi(A \mid S)/B(A \mid S)$

    $V(S) = V(S) + \alpha\rho\delta$

    $\bar{R} = \bar{R} + \eta_r \alpha\rho\delta$

    **for** $i = 1, 2, \ldots, n$ **do**

        $\beta_i = \sum_{j=1}^{m}(-1/b_i^j)(f_j - \bar{R} - \delta)\mathbb{1}\{r_{j-1} \leq R < r_j\}$

        $Z_i = Z_i + \eta_{z_i}\alpha\rho\beta_i$

    **end for**

    $S = S'$

**end while**

return V

---

**Algorithm 2** RED Q-Learning (Tabular)

---

**Input:** the policy $\pi$ to be used (e.g., $\varepsilon$-greedy), piecewise linear subtask function $f$ with $n$ subtasks, $m$ piecewise segments, piecewise conditions $r_{j-1} \leq R < r_j$ such that $f_j$ denotes the $j$th segment of $f$ that satisfies $r_{j-1} \leq R < r_j$, and constants $b_1^j, b_2^j, \ldots, b_n^j \ \forall j = 1, 2, \ldots, m$

**Algorithm parameters:** step size parameters $\alpha, \eta_r, \eta_{z_1}, \eta_{z_2}, \ldots, \eta_{z_n}$

Initialize $Q(s, a) \ \forall s, a; \bar{R}$ arbitrarily (e.g. to zero)

Initialize subtasks $Z_1, Z_2, \ldots, Z_n$ arbitrarily (e.g. to zero)

Obtain initial $S$

**while** still time to train **do**

    $A \leftarrow$ action given by $\pi$ for $S$

    Take action $A$, observe $R, S'$

    $\tilde{R} = f(R, Z_1, Z_2, \ldots, Z_n)$

    $\delta = \tilde{R} - \bar{R} + \max_a Q(S', a) - Q(S, A)$

    $Q(S, A) = Q(S, A) + \alpha\delta$

    $\bar{R} = \bar{R} + \eta_r \alpha\delta$

    **for** $i = 1, 2, \ldots, n$ **do**

        $\beta_i = \sum_{j=1}^{m}(-1/b_i^j)(f_j - \bar{R} - \delta)\mathbb{1}\{r_{j-1} \leq R < r_j\}$

        $Z_i = Z_i + \eta_{z_i}\alpha\beta_i$

    **end for**

    $S = S'$

**end while**

return Q

---

---

**Algorithm 3** RED TD-Learning (Function Approximation)

---

**Input:** the policy $\pi$ to be evaluated, policy $B$ to be used, a differentiable state-value function parameterization: $\hat{v}(s, \boldsymbol{w})$, piecewise linear subtask function $f$ with $n$ subtasks, $m$ piecewise segments, piecewise conditions $r_{j-1} \leq R < r_j$ such that $f_j$ denotes the $j$th segment of $f$ that satisfies $r_{j-1} \leq R < r_j$, and constants $b_1^j, b_2^j, \ldots, b_n^j \ \forall j = 1, 2, \ldots, m$
**Algorithm parameters:** step size parameters $\alpha, \eta_r, \eta_{z_1}, \eta_{z_2}, \ldots, \eta_{z_n}$
Initialize state-value weights $\boldsymbol{w} \in \mathbb{R}^d$ arbitrarily (e.g. to $\boldsymbol{0}$)
Initialize subtasks $Z_1, Z_2, \ldots, Z_n$ arbitrarily (e.g. to zero)
Obtain initial $S$
**while** still time to train **do**
    $A \leftarrow$ action given by $B$ for $S$
    Take action $A$, observe $R, S'$
    $\tilde{R} = f(R, Z_1, Z_2, \ldots, Z_n)$
    $\delta = \tilde{R} - \bar{R} + \hat{v}(S', \boldsymbol{w}) - \hat{v}(S, \boldsymbol{w})$
    $\rho = \pi(A \mid S)/B(A \mid S)$
    $\boldsymbol{w} = \boldsymbol{w} + \alpha \rho \delta \nabla \hat{v}(S, \boldsymbol{w})$
    $\bar{R} = \bar{R} + \eta_r \alpha \rho \delta$
    **for** $i = 1, 2, \ldots, n$ **do**
        $\beta_i = \sum_{j=1}^m (-1/b_i^j)(f_j - \bar{R} - \delta)\mathbb{1}\{r_{j-1} \leq R < r_j\}$
        $Z_i = Z_i + \eta_{z_i} \alpha \rho \beta_i$
    **end for**
    $S = S'$
**end while**
return $\boldsymbol{w}$

---

**Algorithm 4** RED Q-Learning (Function Approximation)

---

**Input:** the policy $\pi$ to be used (e.g., $\varepsilon$-greedy), a differentiable state-action value function parameterization: $\hat{q}(s, a, \boldsymbol{w})$, piecewise linear subtask function $f$ with $n$ subtasks, $m$ piecewise segments, piecewise conditions $r_{j-1} \leq R < r_j$ such that $f_j$ denotes the $j$th segment of $f$ that satisfies $r_{j-1} \leq R < r_j$, and constants $b_1^j, b_2^j, \ldots, b_n^j \ \forall j = 1, 2, \ldots, m$
**Algorithm parameters:** step size parameters $\alpha, \eta_r, \eta_{z_1}, \eta_{z_2}, \ldots, \eta_{z_n}$
Initialize state-action value weights $\boldsymbol{w} \in \mathbb{R}^d$ arbitrarily (e.g. to $\boldsymbol{0}$)
Initialize subtasks $Z_1, Z_2, \ldots, Z_n$ arbitrarily (e.g. to zero)
Obtain initial $S$
**while** still time to train **do**
    $A \leftarrow$ action given by $\pi$ for $S$
    Take action $A$, observe $R, S'$
    $\tilde{R} = f(R, Z_1, Z_2, \ldots, Z_n)$
    $\delta = \tilde{R} - \bar{R} + \max_a \hat{q}(S', a, \boldsymbol{w}) - \hat{q}(S, A, \boldsymbol{w})$
    $\boldsymbol{w} = \boldsymbol{w} + \alpha \delta \nabla \hat{q}(S, A, \boldsymbol{w})$
    $\bar{R} = \bar{R} + \eta_r \alpha \delta$
    **for** $i = 1, 2, \ldots, n$ **do**
        $\beta_i = \sum_{j=1}^m (-1/b_i^j)(f_j - \bar{R} - \delta)\mathbb{1}\{r_{j-1} \leq R < r_j\}$
        $Z_i = Z_i + \eta_{z_i} \alpha \beta_i$
    **end for**
    $S = S'$
**end while**
return $\boldsymbol{w}$

---

## B Convergence Proofs

In this appendix, we present the full convergence proofs for the tabular RED TD-learning and tabular RED Q-learning algorithms. Our general strategy is as follows: we first show that the results from Wan et al. (2021), which show the almost sure convergence of the value function and average-reward estimates of differential algorithms, are applicable to our algorithms. We then build upon these results to show that the subtask estimates of our algorithms converge as well.

For consistency, we adopt similar notation as Wan et al. (2021) for our proofs:

- For a given vector $x$, let $\sum x$ denote the sum of all elements in $x$, such that $\sum x \doteq \sum_i x(i)$.

- Let $\bar{r}_*$ denote the optimal average-reward.

- Let $z_{i_*}$ denote the corresponding optimal subtask value for subtask $z_i \in \mathcal{Z}$.

### B.1 Convergence Proof for the Tabular RED TD-learning Algorithm

In this section, we present the proof for the convergence of the value function, average-reward, and subtask estimates of the RED TD-learning algorithm. Similar to what was done in Wan et al. (2021), we will begin by considering a general algorithm, called *General RED TD*. We will first define General RED TD, then show how the RED TD-learning algorithm is a special case of this algorithm. We will then provide necessary assumptions, state the convergence theorem of General RED TD, and then provide a proof for the theorem, where we show that the value function, average-reward, and subtask estimates converge, thereby showing that the RED TD-learning algorithm converges. We begin by introducing the General RED TD algorithm:

Consider an MDP $\mathcal{M} \doteq \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, p \rangle$, a behavior policy, $B$, and a target policy, $\pi$. Given a state $s \in \mathcal{S}$ and discrete step $n \geq 0$, let $A_n(s) \sim B(\cdot \mid s)$ denote the action selected using the behavior policy, let $R_n(s, A_n(s)) \in \mathcal{R}$ denote a sample of the resulting reward, and let $S'_n(s, A_n(s)) \sim p(\cdot, \cdot \mid s, a)$ denote a sample of the resulting state. Let $\{Y_n\}$ be a set-valued process taking values in the set of nonempty subsets of $\mathcal{S}$, such that: $Y_n = \{s : s$ component of the $|\mathcal{S}|$-sized table of state-value estimates, $V$, that was updated at step $n\}$. Let $\nu(n, s) \doteq \sum_{j=0}^n I\{s \in Y_j\}$, where $I$ is the indicator function, such that $\nu(n, s)$ represents the number of times that $V(s)$ was updated up until step $n$.

Now, let $f$ be a valid subtask function (see Definition 5.1), such that $\tilde{R}_n(s, A_n(s)) \doteq f(R_n(s, A_n(s)), Z_{1,n}, Z_{2,n}, \ldots, Z_{k,n})$ for $k$ subtasks $\in \mathcal{Z}$, where $\tilde{R}_n(s, A_n(s))$ is the extended reward, $\mathcal{Z}$ is the set of subtasks, and $Z_{i,n}$ denotes the estimate of subtask $z_i \in \mathcal{Z}$ at step $n$. Consider an MDP with the extended reward: $\tilde{\mathcal{M}} \doteq \langle \mathcal{S}, \mathcal{A}, \tilde{\mathcal{R}}, \tilde{p} \rangle$, such that $\tilde{R}_n(s, A_n(s)) \in \tilde{\mathcal{R}}$. The update rules of General RED TD for this MDP are as follows, for $n \geq 0$:

$$V_{n+1}(s) \doteq V_n(s) + \alpha_{\nu(n,s)} \rho_n(s) \delta_n(s) I\{s \in Y_n\}, \quad \forall s \in \mathcal{S}, \tag{B.1}$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta_r \sum_s \alpha_{\nu(n,s)} \rho_n(s) \delta_n(s) I\{s \in Y_n\}, \tag{B.2}$$

$$Z_{i,n+1} \doteq Z_{i,n} + \eta_{z_i} \sum_s \alpha_{\nu(n,s)} \rho_n(s) \beta_{i,n}(s) I\{s \in Y_n\}, \quad \forall z_i \in \mathcal{Z}, \tag{B.3}$$

where,

$$\begin{aligned} \delta_n(s) &\doteq \tilde{R}_n(s, A_n(s)) - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s) \\ &= f(R_n(s, A_n(s)), Z_{1,n}, Z_{2,n}, \ldots, Z_{k,n}) - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s), \end{aligned} \tag{B.4}$$

and,

$$\beta_{i,n}(s) \doteq \phi_{i,n}(s) - Z_{i,n}, \quad \forall z_i \in \mathcal{Z}. \tag{B.5}$$

17

536  Here, $\rho_n(s) \doteq \pi(A_n(s) \mid s) \,/\, B(A_n(s) \mid s)$ denotes the importance sampling ratio (with behavior
537  policy, $B$), $\bar{R}_n$ denotes the estimate of the average-reward (see Equation (2)), $\delta_n(s)$ denotes the TD
538  error, $\eta_r$ and $\eta_{z_i}$ are positive scalars, $\phi_{i,n}(s)$ denotes the (potentially-piecewise) subtask target, as
539  defined in Section 5.1, and $\alpha_{\nu(n,s)}$ denotes the step size at time step $n$ for state $s$.

540  We now show that the RED TD-learning algorithm is a special case of the General RED TD algo-
541  rithm. Consider a sequence of experience from our MDP $\tilde{\mathcal{M}}$: $S_t, A_t(S_t), \tilde{R}_{t+1}, S_{t+1}, \dots$ . Now
542  recall the set-valued process $\{Y_n\}$. If we let $n$ = time step $t$, we have:

$$Y_t(s) = \begin{cases} 1, s = S_t, \\ 0, \text{ otherwise,} \end{cases}$$

543  as well as $S'_n(S_t, A_t(S_t)) = S_{t+1}$, $R_n(S_t, A_t) = R_{t+1}$, $\tilde{R}_n(S_t, A_t(S_t)) = \tilde{R}_{t+1}$.
544

545  Hence, update rules (B.1), (B.2), (B.3), (B.4), and (B.5) become:

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha_{\nu(t,S_t)}\rho_t(S_t)\delta_t \text{ , and } V_{t+1}(s) \doteq V_t(s), \forall s \neq S_t, \tag{B.6}$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta_r \alpha_{\nu(t,S_t)}\rho_t(S_t)\delta_t, \tag{B.7}$$
$$Z_{i,t+1} \doteq Z_{i,t} + \eta_{z_i}\alpha_{\nu(t,S_t)}\rho_t(S_t)\beta_{i,t}, \quad \forall z_i \in \mathcal{Z}, \tag{B.8}$$
$$\begin{aligned} \delta_t &\doteq \tilde{R}_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t), \\ &= f(R_{t+1}, Z_{1,t}, Z_{2,t}, \dots, Z_{k,t}) - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t), \end{aligned} \tag{B.9}$$
$$\beta_{i,t} \doteq \phi_{i,t} - Z_{i,t}, \quad \forall z_i \in \mathcal{Z}, \tag{B.10}$$

546  which are RED TD-learning's update rules with $\alpha_{\nu(t,S_t)}$ denoting the step size at time $t$.
547

548  We now specify the assumptions on General RED TD that are needed to ensure convergence. We
549  refer the reader to Wan et al. (2021) for an in-depth discussion on Assumptions B.1 – B.5:
550

551  **Assumption B.1** (Unichain Assumption). *The Markov chain induced by the target policy is*
552  *unichain.*
553

554  **Assumption B.2** (Coverage Assumption). $B(a \mid s) > 0$ *if* $\pi(a \mid s) > 0$ *for all* $s \in \mathcal{S}$, $a \in \mathcal{A}$.
555

556  **Assumption B.3** (Step Size Assumption). $\alpha_n > 0$, $\sum_{n=0}^{\infty} \alpha_n = \infty$, $\sum_{n=0}^{\infty} \alpha_n^2 < \infty$.
557

558  **Assumption B.4** (Asynchronous Step Size Assumption 1). *Let* $[\cdot]$ *denote the integer part of* $(\cdot)$. *For*
559  $x \in (0,1)$,

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

560  *and*

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

561  *uniformly in* $y \in [x, 1]$.
562

563  **Assumption B.5** (Asynchronous Step Size Assumption 2). *There exists* $\Delta > 0$ *such that*

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

564    *a.s., for all $s \in \mathcal{S}$.*

565    *Furthermore, for all $x > 0$, and*

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

566    *the limit*

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

567    *exists a.s. for all $s, s'$.*

568

569    Assumptions B.3, B.4, and B.5, which originate from Borkar (1998), outline the step size require-
570    ments needed to show the convergence of stochastic approximation algorithms. Assumptions B.3
571    and B.4 can be satisfied with step size sequences that decrease to 0 appropriately, including $1/n$,
572    $1/(n \log n)$, and $\log n/n$ (Abounadi et al., 2001). Assumption B.5 first requires that the limiting
573    ratio of visits to any given state, compared to the total number of visits to all states, is greater than or
574    equal to some fixed positive value. The assumption then requires that the relative update frequency
575    between any two states is finite. For instance, Assumption B.5 can be satisfied with $\alpha_n = 1/n$ (see
576    page 403 of Bertsekas and Tsitsiklis (1996) for more information).

577

578    **Assumption B.6** (Subtask Function Assumption). *The subtask function, $f$, is 1) linear or piecewise*
579    *linear, and 2) is invertible with respect to each input given all other inputs.*

580

581    **Assumption B.7** (Subtask Independence Assumption). *Each subtask $z_i \in \mathcal{Z}$ in $f$ is in-*
582    *dependent of the states and actions, and hence independent of the observed reward, $R_n$,*
583    *such that $\tilde{p}(s', f(r, z_1, \ldots, z_n)|s, a) = p(s', r|s, a)$, and $\mathbb{E}[f_j(R_n, Z_{1,n}, Z_{2,n}, \ldots, Z_{k,n})] =$*
584    *$f_j(\mathbb{E}[R_n], Z_{1,n}, Z_{2,n}, \ldots, Z_{k,n})$, where $f_j$ denotes the $j$th segment of a piecewise linear subtask*
585    *function, and $\mathbb{E}$ denotes any expectation taken with respect to the states and actions.*

586

587    **Assumption B.8** (Subtask Step Size Assumption). *If the subtask function is piecewise linear with*
588    *at least two piecewise segments, the subtask step sizes, $\eta_{z_i} \alpha_n$, satisfy the following properties:*
589    *$\eta_{z_i} \alpha_n > 0$, $\sum_{n=0}^{\infty} \eta_{z_i} \alpha_n = \infty$, $\sum_{n=0}^{\infty} (\alpha_n^2 + \eta_{z_i}^2 \alpha_n^2) < \infty$, and $(\eta_{z_i} \alpha_n)/\alpha_n \to 0$, $\forall z_i \in \mathcal{Z}$.*

590    Assumptions B.6, B.7, and B.8 outline the subtask-related requirements. Assumption B.6 ensures
591    that we can explicitly write out the update (B.3), and Assumption B.7 ensures that we do not break
592    the Markov property in the process (i.e., we preserve the Markov property by ensuring that the
593    subtasks are independent of the states and actions, and thereby also independent of the observed
594    reward). Assumption B.8 ensures that the subtask step sizes decrease to 0 appropriately.

595

596    We next point out that it is easy to verify that under Assumption B.1, the following system of
597    equations:

$$\begin{aligned} v_\pi(s) &= \sum_a \pi(a \mid s) \sum_{s', \tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(\tilde{r} - \bar{r}_\pi + v_\pi(s')), \text{ for all } s \in \mathcal{S}, \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a)(f(r, z_1, z_2, \ldots, z_k) - \bar{r}_\pi + v_\pi(s')), \end{aligned} \tag{B.11}$$

598    and,

$$\bar{r}_\pi - \bar{R}_0 = \eta_r \left( \sum v_\pi - \sum V_0 \right), \tag{B.12}$$

$$z_{i,\pi} - Z_{i,0} = \eta_i \left( \sum v_\pi - \sum V_0 \right), \text{ for all } z_i \in \mathcal{Z}, \tag{B.13}$$

599  has a unique solution of $v_\pi$, where $\bar{r}_\pi$ denotes the average-reward induced by following a given
600  policy, $\pi$, and $z_{i,\pi}$ denotes the corresponding subtask value for subtask $z_i \in \mathcal{Z}$. Denote this unique
601  solution of $v_\pi$ as $v_\infty$.
602

603  We are now ready to state the convergence theorem:
604

605  **Theorem B.1.1** (Convergence of General RED TD). *If Assumptions B.1 – B.8 hold, then General*
606  *RED TD (Equations* (B.1) *–* (B.5)*) converges a.s., $\bar{R}_n$ to $\bar{r}_\pi$, $Z_{i,n}$ to $z_{i,\pi}$ $\forall z_i \in \mathcal{Z}$, and $V_n$ to $v_\infty$.*

607  We prove this theorem in the following section. To do so, we first show that General RED TD is of
608  the same form as *General Differential TD* from Wan et al. (2021), thereby allowing us to apply their
609  convergence results for the value function and average-reward estimates of General Differential TD
610  to General RED TD. We then build upon these results, using similar techniques as Wan et al. (2021),
611  to show that the subtask estimates converge as well.

### B.1.1 Proof of Theorem B.1.1 (for *Linear* Subtask Functions)

613  We first provide the proof for *linear* subtask functions, where the the reward-extended TD
614  error can be expressed as a constant, subtask-specific fraction of the regular TD error, such that
615  $\beta_{i,n}(s) = (-1/b_i)\delta_n(s)$. We consider the *piecewise linear* case in Section B.1.2.
616

**Convergence of the average-reward and state-value function estimates:**

618  Consider the increment to $\bar{R}_n$ at each step. We can see from Equation (B.2) that the increment is $\eta_r$
619  times the increment to $V_n$. As such, as was done in Wan et al. (2021), we can write the cumulative
620  increment as follows:

$$\bar{R}_n - \bar{R}_0 = \eta_r \sum_{j=0}^{n-1} \sum_s \alpha_{\nu(j,s)} \rho_j(s) \delta_j(s) I\{s \in Y_j\}$$

$$= \eta_r \left( \sum V_n - \sum V_0 \right)$$

$$\implies \bar{R}_n = \eta_r \sum V_n - \eta_r \sum V_0 + \bar{R}_0 = \eta_r \sum V_n - c_r, \tag{B.14}$$

$$\text{where } c_r \doteq \eta_r \sum V_0 - \bar{R}_0. \tag{B.15}$$

621  Similarly, consider the increment to $Z_{i,n}$ (for an arbitrary subtask $z_i \in \mathcal{Z}$) at each step. As per
622  Remark 5.1, we can write the increment in Equation (B.3) as some constant, subtask-specific fraction
623  of the increment to $V_n$. Consequently, we can write the cumulative increment as follows:

$$Z_{i,n} - Z_{i,0} = \eta_{z_i} \sum_{j=0}^{n-1} \sum_s \alpha_{\nu(j,s)} \rho_j(s) \beta_{i,j}(s) I\{s \in Y_j\}$$

$$= \eta_{z_i} \sum_{j=0}^{n-1} \sum_s \alpha_{\nu(j,s)} \rho_j(s) (-1/b_i) \delta_j(s) I\{s \in Y_j\}$$

$$= \eta_i \left( \sum V_n - \sum V_0 \right)$$

$$\implies Z_{i,n} = \eta_i \sum V_n - \eta_i \sum V_0 + Z_{i,0} = \eta_i \sum V_n - c_i, \tag{B.16}$$

where,

$$c_i \doteq \eta_i \sum V_0 - Z_{i,0}, \text{ and} \tag{B.17}$$

$$\eta_i \doteq (-1/b_i)\eta_{z_i}. \tag{B.18}$$

Now consider the subtask function, $f$. At any given time step, the subtask function can be written as: $f_n = \tilde{R}_n(s, A_n(s)) = b_r R_n(s, A_n(s)) + b_0 + b_1 Z_{1,n} + \ldots + b_k Z_{k,n}$, where $b_r, b_0 \in \mathbb{R}$ and $b_i \in \mathbb{R} \setminus \{0\}$. Given Equation (B.16), we can write the subtask function as follows:

$$f_n = b_r R_n(s, A_n(s)) + b_0 + b_1(\eta_1 \sum V_n - c_1) + \ldots + b_k(\eta_k \sum V_n - c_k)$$

$$= b_r R_n(s, A_n(s)) + \eta_f \sum V_n - c_f, \tag{B.19}$$

where, $\eta_f = \sum_{j=1}^{k} b_j \eta_j$ and $c_f = \sum_{j=1}^{k} b_j c_j - b_0$.

As such, we can substitute $\bar{R}_n$ and $Z_{i,n} \; \forall z_i \in \mathcal{Z}$ in (B.1) with (B.14) and (B.19), respectively, $\forall s \in \mathcal{S}$, which yields:

$$V_{n+1}(s) = V_n(s) + \ldots$$
$$\alpha_{\nu(n,s)}\rho_n(s) \left( b_r R_n(s, A_n(s)) + V_n(S'_n(s, A_n(s))) - V_n(s) - \eta_r \sum V_n + c_r + \eta_f \sum V_n - c_f \right) I\{s \in Y_n\}$$

$$V_{n+1}(s) = V_n(s) + \ldots$$
$$\alpha_{\nu(n,s)}\rho_n(s) \left( b_r R_n(s, A_n(s)) + V_n(S'_n(s, A_n(s))) - V_n(s) - \eta_T \sum V_n + c_T \right) I\{s \in Y_n\}$$

$$V_{n+1}(s) = V_n(s) + \ldots$$
$$\alpha_{\nu(n,s)}\rho_n(s) \left( \hat{R}_n(s, A_n(s)) + V_n(S'_n(s, A_n(s))) - V_n(s) - \eta_T \sum V_n \right) I\{s \in Y_n\}, \tag{B.20}$$

where $\eta_T = \eta_r - \eta_f$, $c_T = c_r - c_f$, and $\hat{R}_n(s, A_n(s)) \doteq b_r R_n(s, A_n(s)) + c_T$.

Equation (B.20) is now in the same form as Equation (B.37) (i.e., General Differential TD) from Wan et al. (2021), who showed that the equation converges a.s. $V_n$ to $v_\infty$ as $n \to \infty$. Moreover, from this result, Wan et al. (2021) showed that $\bar{R}_n$ converges a.s. to $\bar{r}_\pi$ as $n \to \infty$. Given that General RED TD adheres to all the assumptions listed for General Differential TD in Wan et al. (2021), these convergence results apply to General RED TD.

**Convergence of the subtask estimates:**

Let $f(Z_{i,n})$ be shorthand for the subtask function (i.e., $\tilde{R}_n(s, A_n(s))$). We can substitute $Z_{i,n}$ in (B.1) with (B.16) $\forall s \in \mathcal{S}$ as follows:

$$V_{n+1}(s) = V_n(s) + \dots$$
$$\alpha_{\nu(n,s)}\rho_n(s)\left(\tilde{R}_n(s, A_n(s)) - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s)\right) I\{s \in Y_n\}$$

$$\implies V_{n+1}(s) = V_n(s) + \dots$$
$$\alpha_{\nu(n,s)}\rho_n(s)\left(f(Z_{i,n}) - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s)\right) I\{s \in Y_n\}$$

$$\implies V_{n+1}(s) = V_n(s) + \dots$$
$$\alpha_{\nu(n,s)}\rho_n(s)\left(f(\underbrace{\eta_i \sum V_n - c_i}_{\hat{Z}_{i,n}}) - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s)\right) I\{s \in Y_n\}$$

$$\implies V_{n+1}(s) = V_n(s) + \dots$$
$$\alpha_{\nu(n,s)}\rho_n(s)\left(\hat{f}(\hat{Z}_{i,n}) - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s)\right) I\{s \in Y_n\}$$

$$\implies V_{n+1}(s) = V_n(s) + \dots$$
$$\alpha_{\nu(n,s)}\rho_n(s)\left(\hat{R}_n - \bar{R}_n + V_n(S'_n(s, A_n(s))) - V_n(s)\right) I\{s \in Y_n\},$$

$$\tag{B.21}$$

where $\hat{R}_n \doteq \hat{f}(\hat{Z}_{i,n}) = f(Z_{i,n} + c_i) = h(\tilde{R}_n)$. Here, $h(\tilde{R}_n)$ corresponds to the change in $\tilde{R}_n$ due to shifting subtask $Z_{i,n}$ by $c_i$. Denote the inverse of $h(\tilde{R}_n)$ (which exists given Assumption B.6) as $h^{-1}$.

Now consider an MDP, $\hat{\mathcal{M}}$, which has rewards, $\hat{\mathcal{R}}$, corresponding to rewards modified by $h$ from the MDP, $\tilde{\mathcal{M}}$, has the same state and action spaces as $\tilde{\mathcal{M}}$, and has the transition probabilities defined as:

$$\hat{p}(s', h(\hat{r}) \mid s, a) \doteq \tilde{p}(s', \tilde{r} \mid s, a), \tag{B.22}$$

such that $\hat{\mathcal{M}} \doteq \langle \mathcal{S}, \mathcal{A}, \hat{\mathcal{R}}, \hat{p} \rangle$. It is easy to check that the unichain assumption holds for the transformed MDP, $\hat{\mathcal{M}}$. As such, given Assumptions B.6 and B.7, the average-reward induced by following policy $\pi$ for the MDP, $\hat{\mathcal{M}}$, $\hat{\bar{r}}_\pi$, can be written as follows:

$$\hat{\bar{r}}_\pi = h(\bar{r}_\pi). \tag{B.23}$$

Now, because

$$v_\infty(s) = \sum_a \pi(a \mid s) \sum_{s',\tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(\tilde{r} + v_\infty(s') - \bar{r}_\pi) \quad \text{(from (B.11))}$$
$$= \sum_a \pi(a \mid s) \sum_{s',\tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(\tilde{r} + v_\infty(s') - h^{-1}(\hat{\bar{r}}_\pi)) \quad \text{(from (B.23))}$$
$$= \sum_a \pi(a \mid s) \sum_{s',\tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(h(\tilde{r}) + v_\infty(s') - \hat{\bar{r}}_\pi) \quad \text{(by linearity of } h)$$
$$= \sum_a \pi(a \mid s) \sum_{s',\tilde{r}} \hat{p}(s', \tilde{r} \mid s, a)(\tilde{r} + v_\infty(s') - \hat{\bar{r}}_\pi) \quad \text{(from (B.22))},$$

653 we can see that $v_\infty$ is a solution of not just the state-value Bellman equation for the MDP, $\tilde{\mathcal{M}}$, but
654 also the state-value Bellman equation for the transformed MDP, $\hat{\mathcal{M}}$.

655 Next, we can write the subtask value induced by following policy $\pi$ for the MDP, $\hat{\mathcal{M}}$, $\hat{z}_{i,\pi}$, as
656 follows:

$$\hat{z}_{i,\pi} = z_{i,\pi} + c_i. \tag{B.24}$$

657 We can then combine Equations (B.13), (B.16), and (B.24), which yields:

$$\hat{z}_{i,\pi} = \eta_i \sum v_\infty. \tag{B.25}$$

658 Next, we can combine Equation (B.16) with the result from Wan et al. (2021) which shows that
659 $V_n \to v_\infty$, which yields:

$$Z_{i,n} \to \eta_i \sum v_\infty - c_i. \tag{B.26}$$

660 Moreover, because $\hat{z}_{i,\pi} = \eta_i \sum v_\infty$ (Equation (B.25)), we have:

$$Z_{i,n} \to \hat{z}_{i,\pi} - c_i. \tag{B.27}$$

661 Finally, because $\hat{z}_{i,\pi} = z_{i,\pi} + c_i$ (Equation (B.24)), we have:

$$Z_{i,n} \to z_{i,\pi} \quad \text{a.s. as} \quad n \to \infty. \tag{B.28}$$

### B.1.2 Proof of Theorem B.1.1 (for *Piecewise Linear* Subtask Functions)

We now provide the proof for *piecewise linear* subtask functions, where the the reward-extended TD error can be expressed as follows: $\beta_{i,n}(s) = (-1/b_{i,n})(\tilde{R}_n(s, A_n(s)) - \bar{R}_n - \delta_n(s))$. Our general strategy in this case is to use a two-timescales argument, such that we leverage Theorem 2 in Section 6 of Borkar (2009), along with the results from Theorem B.3 of Wan et al. (2021).

To begin, let us consider Assumption B.8. In particular, $(\eta_{z_i}\alpha_n)/\alpha_n \to 0$ implies that the subtask step sizes, $\eta_{z_i}\alpha_n$, decrease to 0 at a faster rate than the value function step size, $\alpha_n$. This implies that the subtask updates move on a slower timescale compared to the value function update. Hence, as argued in Section 6 of Borkar (2009), the (faster) value function update (B.1) views the (slower) subtask updates (B.3) as quasi-static, while the (slower) subtask updates view the (faster) value function update as nearly equilibrated (as we will show below, the results from Wan et al. (2021) imply the existence of such an equilibrium point).

**Convergence of the average-reward and state-action value function estimates:**

Given the two-timescales argument, Equation (B.1) can be viewed as being of the same form as Equation (B.30) (i.e., General Differential TD) from Wan et al. (2021), who showed that the equation converges a.s. $V_n$ to $v_\infty$ as $n \to \infty$. Moreover, from this result, Wan et al. (2021) showed that $\bar{R}_n$ converges a.s. to $\bar{r}_\pi$ as $n \to \infty$. Given that General RED TD adheres to all the assumptions listed for General Differential TD in Wan et al. (2021), these convergence results apply to General RED TD.

**Convergence of the subtask estimates:**

Let us consider the asynchronous subtask updates (B.3). These updates are (each) of the same form as Equation 7.1.2 of Borkar (2009). As such, to show the convergence of the subtask estimates, we can apply the result in Section 7.4 of Borkar (2009), which shows the convergence of asynchronous updates of the same form as Equation 7.1.2. To apply this result, given Assumptions B.4 and B.5, we only need to show the convergence of the *synchronous* version of the subtask updates:

$$Z_{i,n+1} = Z_{i,n} + \eta_{z_i}\alpha_n \left[ (-1/b_{i,n}) \left( \rho_n(\tilde{R}_n - \bar{R}_n) - (g(V_n) + M_{n+1}) \right) \right] \; \forall z_i \in \mathcal{Z} \quad (B.29)$$

where,

$$g(V_n)(s) \doteq \sum_{s',\tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(\tilde{r} + V_n(s')) - V_n(s) - \bar{R}_n$$

$$= T(V_n)(s) - V_n(s) - \bar{R}_n, \text{ and}$$

$$M_{n+1}(s) \doteq \rho_n(s) \left( \tilde{R}_n(s, A_n(s)) + V_n(S'_n(s, A_n(s))) - V_n(s) - \bar{R}_n \right) - g(V_n)(s).$$

To show the convergence of the synchronous update (B.29) under the two-timescales argument, we can apply the result of Theorem 2 in Section 6 of Borkar (2009) to show that $Z_{i,n} \to z_{i,\pi} \forall z_i \in \mathcal{Z}$ a.s. as $n \to \infty$. This theorem requires that 3 assumptions be satisfied. As such, we will now show, via Lemmas B.1 - B.3, that these 3 assumptions are indeed satisfied.

**Lemma B.1.** *The value function update rule, $V_{n+1} = V_n + \alpha_n(g(V_n) + M_{n+1})$, has a globally asymptotically stable equilibrium, $v_\infty$.*

*Proof.* This was shown in Theorem B.3 of Wan et al. (2021). □

**Lemma B.2.** *The subtask update rules (B.29) each have a globally asymptotically stable equilibrium, $z_{i,\pi}$.*

*Proof.* Applying the results of Theorem B.3 of Wan et al. (2021) under the two-timescales argument, we have that $g(V_n) \to 0$, $\bar{R}_n \to \bar{r}_\pi$, that $\{M_n\}$ is a martingale difference sequence, such that $\mathbb{E}[M_{n+1} \mid \mathcal{F}_n] = 0$ a.s., $n \geq 0$, and that $\{M_n\}$ is square-integrable, such that $\mathbb{E}[||M_{n+1}||^2 \mid \mathcal{F}_n] \leq K(1 + ||Q_n||^2)$ a.s., $n \geq 0$, for some constant $K > 0$. Given these results, the remaining $\rho_n(s)(\tilde{R}_n(s, A_n(s)) - \bar{R}_n) = \rho_n(s)\tilde{R}_n(s, A_n(s)) - \bar{r}_\pi$ term in the subtask updates (B.29) can be interpreted as a martingale difference sequence, $\{M_n^r\}$, such that $\mathbb{E}[M_{n+1}^r \mid \mathcal{F}_n] = \mathbb{E}[\rho_n(s)(\tilde{R}_{n+1}(s, A_n(s)) - \bar{R}_n) \mid \mathcal{F}_n] = \mathbb{E}[\rho_n(s)\tilde{R}_{n+1}(s, A_n(s)) \mid \mathcal{F}_n] - \bar{r}_\pi = 0$ a.s., $n \geq 0$. As such, given Assumptions B.4, B.5, and B.8, to show that the subtask update rules (B.29) each have a globally asymptotically stable equilibrium, we only need to show that the martingale difference sequence, $\{M_n^r\}$, is square-integrable, such that $\mathbb{E}[(M_{n+1}^r)^2 \mid \mathcal{F}_n] < \infty$ a.s., $n \geq 0$. Indeed, because $\tilde{R}_n(s, A_n(s))$ is bounded, it directly follows that its mean, $\bar{r}_\pi$, is also bounded, and as such, we have that the martingale difference sequence, $\{M_n^r\}$, is square-integrable. Hence, we can conclude that the subtask update rules (B.29) each have a globally asymptotically stable equilibrium, $z_{i,\pi}$. $\qquad\square$

**Lemma B.3.** $\sup_n(||V_n|| + ||Z_n||) < \infty$ a.s.

*Proof.* It was shown in Theorem B.3 of Wan et al. (2021) that $\sup_n(||V_n||) < \infty$ a.s. Hence, we only need to show that $\sup_n(||Z_n||) < \infty$ a.s. To this end, we can apply Theorem 7 in Section 3 of Borkar (2009). This theorem requires 4 assumptions:

- **(A1)** The function $g$ is Lipschitz: $||g(x) - g(y)|| \leq L||x - y||$ for some $0 < L < \infty$.

- **(A2)** The sequence $\{\eta_{z_i}\alpha_n\}$ satisfies $\eta_{z_i}\alpha_n > 0$, $\sum \eta_{z_i}\alpha_n = \infty$, and $\sum \eta_{z_i}^2 \alpha_n^2 < \infty$.

- **(A3)** $\{M_n\}$ and $\{M_n^r\}$ are martingale difference sequences that are square-integrable.

- **(A4)** The functions $g_d(x) \doteq g(dx)/d$, $d \geq 1, x \in \mathbb{R}^k$, satisfy $g_d(x) \to g_\infty(x)$ as $d \to \infty$, uniformly on compacts for some $g_\infty \in C(\mathbb{R}^k)$. Furthermore, the ODE $\dot{x}_t = g_\infty(x_t)$ has the origin as its unique globally asymptotically stable equilibrium.

Under the two-timescales argument, the results of Theorem B.3 of Wan et al. (2021) apply, thereby satisfying the above assumptions, except for the assumptions regarding $\{\eta_{z_i}\alpha_n\}$ and $\{M_n^r\}$. In this regard, Assumptions B.4 and B.8 satisfy Assumption **(A2)**. Moreover, we showed in Lemma B.2 that $\{M_n^r\}$ is indeed a martingale difference sequence that is square-integrable. As such, Assumptions **(A1)** - **(A4)** are verified, meaning that we can apply the results of Theorem 7 in Section 3 of Borkar (2009) to conclude that $\sup_n(||Z_n||) < \infty$ a.s., and hence, that $\sup_n(||V_n|| + ||Z_n||) < \infty$ a.s. $\qquad\square$

As such, we have now verified the 3 assumptions required by Theorem 2 in Section 6 of Borkar (2009), which means that we can apply the result of the theorem to conclude that $Z_{i,n} \to z_{i,\pi} \forall z_i \in \mathcal{Z}$ a.s. as $n \to \infty$.

This completes the proof of Theorem B.1.1.

**B.2 Convergence Proof for the Tabular RED Q-learning Algorithm**

In this section, we present the proof for the convergence of the value function, average-reward, and subtask estimates of the RED Q-learning algorithm. Similar to what was done in Wan et al. (2021), we will begin by considering a general algorithm, called *General RED Q*. We will first define General RED Q, then show how the RED Q-learning algorithm is a special case of this algorithm. We will then provide necessary assumptions, state the convergence theorem of General RED Q, and then provide a proof for the theorem, where we show that the value function, average-reward, and subtask estimates converge, thereby showing that the RED Q-learning algorithm converges. We begin by introducing the General RED Q algorithm:

Consider an MDP $\mathcal{M} \doteq \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, p \rangle$. Given a state $s \in \mathcal{S}$, action $a \in \mathcal{A}$, and discrete step $n \geq 0$, let $R_n(s, a) \in \mathcal{R}$ denote a sample of the resulting reward, and let $S'_n(s, a) \sim p(\cdot \mid s, a)$ denote a sample of the resulting state. Let $\{Y_n\}$ be a set-valued process taking values in the set of nonempty subsets of $\mathcal{S} \times \mathcal{A}$, such that: $Y_n = \{(s, a) : (s, a) \text{ component of the } |\mathcal{S} \times \mathcal{A}|\text{-sized table of state-action value estimates, } Q, \text{ that was updated at step } n\}$. Let $\nu(n, s, a) \doteq \sum_{j=0}^{n} I\{(s, a) \in Y_j\}$, where $I$ is the indicator function, such that $\nu(n, s, a)$ represents the number of times that the $(s, a)$ component of $Q$ was updated up until step $n$.

Now, let $f$ be a valid subtask function (see Definition 5.1), such that $\tilde{R}_n(s, a) \doteq f(R_n(s, a), Z_{1,n}, Z_{2,n}, \ldots, Z_{n,k})$ for $k$ subtasks $\in \mathcal{Z}$, where $\tilde{R}_n(s, a)$ is the extended reward, $\mathcal{Z}$ is the set of subtasks, and $Z_{i,n}$ denotes the estimate of subtask $z_i \in \mathcal{Z}$ at step $n$. Consider an MDP with the extended reward: $\tilde{\mathcal{M}} \doteq \langle \mathcal{S}, \mathcal{A}, \tilde{\mathcal{R}}, \tilde{p} \rangle$, such that $\tilde{R}_n(s, a) \in \tilde{\mathcal{R}}$. The update rules of General RED Q for this MDP are as follows, for $n \geq 0$:

$$Q_{n+1}(s, a) \doteq Q_n(s, a) + \alpha_{\nu(n,s,a)} \delta_n(s, a) I\{(s, a) \in Y_n\}, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}, \tag{B.30}$$

$$\bar{R}_{n+1} \doteq \bar{R}_n + \eta_r \sum_{s,a} \alpha_{\nu(n,s,a)} \delta_n(s, a) I\{(s, a) \in Y_n\}, \tag{B.31}$$

$$Z_{i,n+1} \doteq Z_{i,n} + \eta_{z_i} \sum_{s,a} \alpha_{\nu(n,s,a)} \beta_{i,n}(s, a) I\{(s, a) \in Y_n\}, \quad \forall z_i \in \mathcal{Z} \tag{B.32}$$

where,

$$\begin{aligned} \delta_n(s, a) &\doteq \tilde{R}_n(s, a) - \bar{R}_n + \max_{a'} Q_n(S'_n(s, a), a') - Q_n(s, a) \\ &= f(R_n(s, a), Z_{1,n}, Z_{2,n}, \ldots, Z_{k,n}) - \bar{R}_n + \max_{a'} Q_n(S'_n(s, a), a') - Q_n(s, a), \end{aligned} \tag{B.33}$$

and,

$$\beta_{i,n}(s, a) \doteq \phi_{i,n}(s, a) - Z_{i,n}, \quad \forall z_i \in \mathcal{Z}. \tag{B.34}$$

Here, $\bar{R}_n$ denotes the estimate of the average-reward (see Equation (2)), $\delta_n(s, a)$ denotes the TD error, $\eta_r$ and $\eta_{z_i}$ are positive scalars, $\phi_{i,n}(s, a)$ denotes the (potentially-piecewise) subtask target, as defined in Section 5.1, and $\alpha_{\nu(n,s,a)}$ denotes the step size at time step $n$ for state-action pair $(s, a)$.

We now show that the RED Q-learning algorithm is a special case of the General RED Q algorithm. Consider a sequence of experience from our MDP $\tilde{\mathcal{M}}$: $S_t, A_t, \tilde{R}_{t+1}, S_{t+1}, \ldots$. Now recall the set-valued process $\{Y_n\}$. If we let $n$ = time step $t$, we have:

$$Y_t(s, a) = \begin{cases} 1, & s = S_t \text{ and } a = A_t, \\ 0, & \text{otherwise}, \end{cases}$$

as well as $S'_n(S_t, A_t) = S_{t+1}$, $R_n(S_t, A_t) = R_{t+1}$, and $\tilde{R}_n(S_t, A_t) = \tilde{R}_{t+1}$.

768    Hence, update rules (B.30), (B.31), (B.32), (B.33), and (B.34) become:

$$Q_{t+1}(S_t, A_t) \doteq Q_t(S_t, A_t) + \alpha_{\nu(t, S_t, A_t)} \delta_t; \quad Q_{t+1}(s, a) \doteq Q_t(s, a), \forall s \neq S_t, a \neq A_t, \quad \text{(B.35)}$$

$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta_r \alpha_{\nu(t, S_t, A_t)} \delta_t, \quad \text{(B.36)}$$

$$Z_{i,t+1} \doteq Z_{i,t} + \eta_{z_i} \alpha_{\nu(t, S_t, A_t)} \beta_{i,t}, \quad \forall z_i \in \mathcal{Z}, \quad \text{(B.37)}$$

$$\delta_t \doteq \tilde{R}_{t+1} - \bar{R}_t + \max_{a'} Q_t(S_{t+1}, a') - Q_t(S_t, A_t),$$

$$= f(R_{t+1}, Z_{1,t}, Z_{2,t}, \dots, Z_{k,t}) - \bar{R}_t + \max_{a'} Q_t(S_{t+1}, a') - Q_t(S_t, A_t), \quad \text{(B.38)}$$

$$\beta_{i,t} \doteq \phi_{i,t} - Z_{i,t}, \quad \forall z_i \in \mathcal{Z}, \quad \text{(B.39)}$$

769    which are RED Q-learning's update rules with $\alpha_{\nu(t, S_t, A_t)}$ denoting the step size at time $t$.
770

771    We now specify the assumptions on General RED Q that are needed to ensure convergence. We
772    refer the reader to Wan et al. (2021) for an in-depth discussion on these assumptions:
773

774    **Assumption B.9** (Communicating Assumption). *The MDP has a single communicating class. That*
775 *is, each state in the MDP is accessible from every other state under some deterministic stationary*
776 *policy.*
777

778    **Assumption B.10** (State-Action Value Function Uniqueness). *There exists a unique solution of $q$*
779 *only up to a constant in the Bellman equation (4).*
780

781    **Assumption B.11** (Asynchronous Step Size Assumption 3). *There exists $\Delta > 0$ such that*

$$\liminf_{n \to \infty} \frac{\nu(n, s, a)}{n + 1} \geq \Delta,$$

782 *a.s., for all $s \in \mathcal{S}, a \in \mathcal{A}$.*
783

784 *Furthermore, for all $x > 0$, and*

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

785 *the limit*

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

786 *exists a.s. for all $s, s', a, a'$.*
787

788    We next point out that it is easy to verify that under Assumption B.9, the following system of
789 equations:

$$q_\pi(s, a) = \sum_{s', \tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(\tilde{r} - \bar{r}_\pi + \max_{a'} q_\pi(s, a)), \quad \forall s \in \mathcal{S}, a \in \mathcal{A},$$

$$= \sum_{s', r} p(s', r \mid s, a)(f(r, z_1, z_2, \dots, z_k) - \bar{r}_\pi + \max_{a'} q_\pi(s, a)), \quad \text{(B.40)}$$

790 and,

$$\bar{r}_* - \bar{R}_0 = \eta_r \left( \sum q_\pi - \sum Q_0 \right), \tag{B.41}$$

$$z_{i_*} - Z_{i,0} = \eta_i \left( \sum q_\pi - \sum Q_0 \right), \quad \forall z_i \in \mathcal{Z}, \tag{B.42}$$

791 has a unique solution for $q_\pi$, where $\bar{r}_*$ denotes the optimal average-reward, and $z_{i_*}$ denotes the
792 corresponding optimal subtask value for subtask $z_i \in \mathcal{Z}$. Denote this unique solution for $q_\pi$ as $q_*$.
793

794 We are now ready to state the convergence theorem:
795

796 **Theorem B.2.1** (Convergence of General RED Q). *If Assumptions B.3, B.4, B.6, B.7, B.8, B.9, B.10,*
797 *and B.11 hold, then the General RED Q algorithm (Equations B.30–B.34) converges a.s. $\bar{R}_n$ to $\bar{r}_*$,*
798 *$Z_{i,n}$ to $z_{i_*}$ $\forall z_i \in \mathcal{Z}$, $Q_n$ to $q_*$, $\bar{r}_{\pi_t}$ to $\bar{r}_*$, and $z_{i,\pi_t}$ to $z_{i_*}$ $\forall z_i \in \mathcal{Z}$, where $\pi_t$ is any greedy policy*
799 *with respect to $Q_t$, and $z_{i,\pi_t}$ denotes the subtask value induced by following policy $\pi_t$.*

800 We prove this theorem in the following section. To do so, we first show that General RED Q is of
801 the same form as *General Differential Q* from Wan et al. (2021), thereby allowing us to apply their
802 convergence results for the value function and average-reward estimates of General Differential Q
803 to General RED Q. We then build upon these results, using similar techniques as Wan et al. (2021),
804 to show that the subtask estimates converge as well.

### B.2.1 Proof of Theorem B.2.1 (for *Linear* Subtask Functions)

806 We first provide the proof for *linear* subtask functions, where the the reward-extended TD
807 error can be expressed as a constant, subtask-specific fraction of the regular TD error, such that
808 $\beta_{i,n}(s,a) = (-1/b_i)\delta_n(s,a)$. We consider the *piecewise linear* case in Section B.2.2.
809

**Convergence of the average-reward and state-action value function estimates:**

811 Consider the increment to $\bar{R}_n$ at each step. We can see from Equation (B.31) that the increment is $\eta_r$
812 times the increment to $Q_n$. As such, as was done in Wan et al. (2021), we can write the cumulative
813 increment as follows:

$$\bar{R}_n - \bar{R}_0 = \eta_r \sum_{j=0}^{n-1} \sum_{s,a} \alpha_{\nu(j,s,a)} \delta_j(s,a) I\{(s,a) \in Y_j\}$$

$$= \eta_r \left( \sum Q_n - \sum Q_0 \right)$$

$$\implies \bar{R}_n = \eta_r \sum Q_n - \eta_r \sum Q_0 + \bar{R}_0 = \eta_r \sum Q_n - c_r, \tag{B.43}$$

$$\text{where } c_r \doteq \eta_r \sum Q_0 - \bar{R}_0. \tag{B.44}$$

814 Similarly, consider the increment to $Z_{i,n}$ (for an arbitrary subtask $z_i \in \mathcal{Z}$) at each step. As per Re-
815 mark 5.1, we can write the increment in Equation (B.32) as some constant, subtask-specific fraction
816 of the increment to $Q_n$. Consequently, we can write the cumulative increment as follows:

$$Z_{i,n} - Z_{i,0} = \eta_{z_i} \sum_{j=0}^{n-1} \sum_{s,a} \alpha_{\nu(j,s,a)} \beta_{i,j}(s,a) I\{(s,a) \in Y_j\}$$

$$= \eta_{z_i} \sum_{j=0}^{n-1} \sum_{s,a} \alpha_{\nu(j,s,a)} (-1/b_i) \delta_j(s,a) I\{(s,a) \in Y_j\}$$

$$= \eta_i \left( \sum Q_n - \sum Q_0 \right)$$

$$\implies Z_{i,n} = \eta_i \sum Q_n - \eta_i \sum Q_0 + Z_{i,0} = \eta_i \sum Q_n - c_i, \tag{B.45}$$

where,

$$c_i \doteq \eta_i \sum Q_0 - Z_{i,0}, \text{ and} \tag{B.46}$$

$$\eta_i \doteq (-1/b_i)\eta_{z_i}. \tag{B.47}$$

Now consider the subtask function, $f$. At any given time step, the subtask function can be written as: $f_n = \tilde{R}_n(s,a) = b_r R_n(s,a) + b_0 + b_1 Z_{1,n} + \ldots + b_k Z_{k,n}$, where $b_r, b_0 \in \mathbb{R}$ and $b_i \in \mathbb{R} \setminus \{0\}$. Given Equation (B.45), we can write the subtask function as follows:

$$f_n = b_r R_n(s,a) + b_0 + b_1(\eta_1 \sum Q_n - c_1) + \ldots + b_k(\eta_k \sum Q_n - c_k)$$
$$= b_r R_n(s,a) + \eta_f \sum Q_n - c_f, \tag{B.48}$$

where, $\eta_f = \sum_{j=1}^{k} b_j \eta_j$ and $c_f = \sum_{j=1}^{k} b_j c_j - b_0$.

As such, we can substitute $\bar{R}_n$ and $Z_{i,n}$ $\forall z_i \in \mathcal{Z}$ in (B.30) with (B.43) and (B.48), respectively, $\forall s \in \mathcal{S}, a \in \mathcal{A}$, which yields:

$$Q_{n+1}(s,a) = Q_n(s,a) + \ldots$$
$$\alpha_{\nu(n,s,a)} \left( b_r R_n(s,a) + \max_{a'} Q_n(S'_n(s,a), a') - Q_n(s,a) - \eta_r \sum Q_n + c_r + \eta_f \sum Q_n - c_f \right) I\{(s,a) \in Y_n\}$$

$$Q_{n+1}(s,a) = Q_n(s,a) + \ldots$$
$$\alpha_{\nu(n,s,a)} \left( b_r R_n(s,a) + \max_{a'} Q_n(S'_n(s,a), a') - Q_n(s,a) - \eta_T \sum Q_n + c_T \right) I\{(s,a) \in Y_n\}$$

$$Q_{n+1}(s,a) = Q_n(s,a) + \ldots$$
$$\alpha_{\nu(n,s,a)} \left( \hat{R}_n(s,a) + \max_{a'} Q_n(S'_n(s,a), a') - Q_n(s,a) - \eta_T \sum Q_n \right) I\{(s,a) \in Y_n\}, \tag{B.49}$$

where $\eta_T = \eta_r - \eta_f$, $c_T = c_r - c_f$, and $\hat{R}_n(s,a) \doteq b_r R_n(s,a) + c_T$.

826    Equation (B.49) is now in the same form as Equation (B.14) (i.e., General Differential Q) from Wan
827    et al. (2021), who showed that the equation converges a.s. $Q_n$ to $q_*$ as $n \to \infty$. Moreover, from this
828    result, Wan et al. (2021) showed that $\bar{R}_n$ converges a.s. to $\bar{r}_*$ as $n \to \infty$, and that $\bar{r}_{\pi_t}$ converges a.s.
829    to $\bar{r}_*$, where $\pi_t$ is a greedy policy with respect to $Q_t$. Given that General RED Q adheres to all the
830    assumptions listed for General Differential Q in Wan et al. (2021), these convergence results apply
831    to General RED Q.
832

**Convergence of the subtask estimates:**

834    Let $f(Z_{i,n})$ be shorthand for the subtask function (i.e., $\tilde{R}_n(s,a)$). We can substitute $Z_{i,n}$ in (B.30)
835    with (B.45) $\forall s \in \mathcal{S}, a \in \mathcal{A}$ as follows:

$$Q_{n+1}(s,a) = Q_n(s,a) + \dots$$
$$\alpha_{\nu(n,s,a)} \left( \tilde{R}_n(s,a) - \bar{R}_n + \max_{a'} Q_n(S_n'(s,a),a') - Q_n(s,a) \right) I\{(s,a) \in Y_n\}$$

$$\implies Q_{n+1}(s,a) = Q_n(s,a) + \dots$$
$$\alpha_{\nu(n,s,a)} \left( f(Z_{i,n}) - \bar{R}_n + \max_{a'} Q_n(S_n'(s,a),a') - Q_n(s,a) \right) I\{(s,a) \in Y_n\}$$

$$\implies Q_{n+1}(s,a) = Q_n(s,a) + \dots$$
$$\alpha_{\nu(n,s,a)} \left( f(\underbrace{\eta_i \sum Q_n}_{\hat{Z}_{i,n}} -c_i) - \bar{R}_n + \max_{a'} Q_n(S_n'(s,a),a') - Q_n(s,a) \right) I\{(s,a) \in Y_n\}$$

$$\implies Q_{n+1}(s,a) = Q_n(s,a) + \dots$$
$$\alpha_{\nu(n,s,a)} \left( \hat{f}(\hat{Z}_{i,n}) - \bar{R}_n + \max_{a'} Q_n(S_n'(s,a),a') - Q_n(s,a) \right) I\{(s,a) \in Y_n\}$$

$$\implies Q_{n+1}(s,a) = Q_n(s,a) + \dots$$
$$\alpha_{\nu(n,s,a)} \left( \hat{R}_n - \bar{R}_n + \max_{a'} Q_n(S_n'(s,a),a') - Q_n(s,a) \right) I\{(s,a) \in Y_n\},$$

$$\tag{B.50}$$

836    where $\hat{R}_n \doteq \hat{f}(\hat{Z}_{i,n}) = f(Z_{i,n} + c_i) = h(\tilde{R}_n)$. Here, $h(\tilde{R}_n)$ corresponds to the change in $\tilde{R}_n$ due
837    to shifting subtask $Z_{i,n}$ by $c_i$. Denote the inverse of $h(\tilde{R}_n)$ (which exists given Assumption B.6) as
838    $h^{-1}$.
839

840    Now consider an MDP, $\hat{\mathcal{M}}$, which has rewards, $\hat{\mathcal{R}}$, corresponding to rewards modified by $h$ from the
841    MDP, $\tilde{\mathcal{M}}$, has the same state and action spaces as $\tilde{\mathcal{M}}$, and has the transition probabilities defined as:

$$\hat{p}\left(s', h(\tilde{r}) \mid s,a\right) \doteq \tilde{p}(s', \tilde{r} \mid s,a), \tag{B.51}$$

842    such that $\hat{\mathcal{M}} \doteq \langle \mathcal{S}, \mathcal{A}, \hat{\mathcal{R}}, \hat{p} \rangle$. It is easy to check that the communicating assumption holds for the
843    transformed MDP, $\hat{\mathcal{M}}$. As such, given Assumptions B.6 and B.7, the optimal average-reward for the
844    MDP, $\hat{\mathcal{M}}, \hat{\bar{r}}_*$, can be written as follows:

$$\hat{\bar{r}}_* = h(\bar{r}_*). \tag{B.52}$$

Now, because

$$q_*(s,a) = \sum_{s',\tilde{r}} \tilde{p}(s',\tilde{r} \mid s,a)(\tilde{r} + \max_{a'} q_*(s',a') - \bar{r}_*) \quad \text{(from (B.40))}$$

$$= \sum_{s',\tilde{r}} \tilde{p}(s',\tilde{r} \mid s,a)(\tilde{r} + \max_{a'} q_*(s',a') - h^{-1}(\hat{\bar{r}}_*)) \quad \text{(from (B.52))}$$

$$= \sum_{s',\tilde{r}} \tilde{p}(s',\tilde{r} \mid s,a)(h(\tilde{r}) + \max_{a'} q_*(s',a') - \hat{\bar{r}}_*) \quad \text{(by linearity of } h)$$

$$= \sum_{s',\tilde{r}} \hat{p}(s',\tilde{r} \mid s,a)(\tilde{r} + \max_{a'} q_*(s',a') - \hat{\bar{r}}_*) \quad \text{(from (B.51))},$$

we can see that $q_*$ is a solution of not just the state-action value Bellman equation for the MDP, $\tilde{\mathcal{M}}$, but also the state-action value Bellman equation for the transformed MDP, $\hat{\mathcal{M}}$.

Next, we can write the optimal subtask value for the MDP, $\hat{\mathcal{M}}$, $\hat{z}_{i_*}$, as follows:

$$\hat{z}_{i_*} = z_{i_*} + c_i. \tag{B.53}$$

We can then combine Equations (B.42), (B.45), and (B.53), which yields:

$$\hat{z}_{i_*} = \eta_i \sum q_*. \tag{B.54}$$

Next, we can combine Equation (B.45) with the result from Wan et al. (2021) which shows that $Q_n \to q_*$, which yields:

$$Z_{i,n} \to \eta_i \sum q_* - c_i. \tag{B.55}$$

Moreover, because $\eta_i \sum q_* = \hat{z}_{i_*}$ (Equation (B.54)), we have:

$$Z_{i,n} \to \hat{z}_{i_*} - c_i. \tag{B.56}$$

Finally, because $\hat{z}_{i_*} = z_{i_*} + c_i$ (Equation (B.53)), we have:

$$Z_{i,n} \to z_{i_*} \quad \text{a.s. as} \quad n \to \infty. \tag{B.57}$$

We conclude by considering $z_{i,\pi_t} \; \forall z_i \in \mathcal{Z}$, where $\pi_t$ is a greedy policy with respect to $Q_t$. Given that $Q_t \to q_*$ and $\bar{r}_{\pi_t} \to \bar{r}_*$ a.s., it directly follows from Definition 5.1 that $z_{i,\pi_t} \to z_{i_*} \; \forall z_i \in \mathcal{Z}$ a.s.

### B.2.2 Proof of Theorem B.2.1 (for *Piecewise Linear* Subtask Functions)

We now provide the proof for *piecewise linear* subtask functions, where the the reward-extended TD error can be expressed as follows: $\beta_{i,n}(s,a) = (-1/b_{i,n})(\tilde{R}_n(s,a) - \bar{R}_n - \delta_n(s,a))$. Our general strategy in this case is to use a two-timescales argument, such that we leverage Theorem 2 in Section 6 of Borkar (2009), along with the results from Theorems B.1 and B.2 of Wan et al. (2021).

To begin, let us consider Assumption B.8. In particular, $(\eta_{z_i}\alpha_n)/\alpha_n \to 0$ implies that the subtask step sizes, $\eta_{z_i}\alpha_n$, decrease to 0 at a faster rate than the value function step size, $\alpha_n$. This implies that the subtask updates move on a slower timescale compared to the value function update. Hence, as argued in Section 6 of Borkar (2009), the (faster) value function update (B.30) views the (slower) subtask updates (B.32) as quasi-static, while the (slower) subtask updates view the (faster) value function update as nearly equilibrated (as we will show below, the results from Wan et al. (2021) imply the existence of such an equilibrium point).

870 **Convergence of the average-reward and state-action value function estimates:**

871 Given the two-timescales argument, Equation (B.30) can be viewed as being of the same form as
872 Equation (B.4) (i.e., General Differential Q) from Wan et al. (2021), who showed that the equation
873 converges a.s. $Q_n$ to $q_*$ as $n \to \infty$. Moreover, from this result, Wan et al. (2021) showed that
874 $\bar{R}_n$ converges a.s. to $\bar{r}_*$ as $n \to \infty$, and that $\bar{r}_{\pi_t}$ converges a.s. to $\bar{r}_*$, where $\pi_t$ is a greedy policy
875 with respect to $Q_t$. Given that General RED Q adheres to all the assumptions listed for General
876 Differential Q in Wan et al. (2021), these convergence results apply to General RED Q.

877

878 **Convergence of the subtask estimates:**

879 Let us consider the asynchronous subtask updates (B.32). These updates are (each) of the same form
880 as Equation 7.1.2 of Borkar (2009). As such, to show the convergence of the subtask estimates, we
881 can apply the result in Section 7.4 of Borkar (2009), which shows the convergence of asynchronous
882 updates of the same form as Equation 7.1.2. To apply this result, given Assumptions B.4 and B.11,
883 we only need to show the convergence of the *synchronous* version of the subtask updates:

$$Z_{i,n+1} = Z_{i,n} + \eta_{z_i}\alpha_n \left[ (-1/b_{i,n}) \left( \tilde{R}_n - \bar{R}_n - (g(Q_n) + M_{n+1}) \right) \right] \ \forall z_i \in \mathcal{Z} \qquad \text{(B.58)}$$

884 where,

$$g(Q_n)(s,a) \doteq \sum_{s',\tilde{r}} \tilde{p}(s', \tilde{r} \mid s, a)(\tilde{r} + \max_{a'} Q_n(s', a')) - Q_n(s,a) - \bar{R}_n$$

$$= T(Q_n)(s,a) - Q_n(s,a) - \bar{R}_n, \text{ and}$$

$$M_{n+1}(s,a) \doteq \tilde{R}_n(s,a) + \max_{a'} Q_n(S'_n(s,a), a') - T(Q_n)(s,a).$$

885 To show the convergence of the synchronous update (B.58) under the two-timescales argument, we
886 can apply the result of Theorem 2 in Section 6 of Borkar (2009) to show that $Z_{i,n} \to z_{i_*} \forall z_i \in \mathcal{Z}$
887 a.s. as $n \to \infty$. This theorem requires that 3 assumptions be satisfied. As such, we will now show,
888 via Lemmas B.4 - B.6, that these 3 assumptions are indeed satisfied.

889

890 **Lemma B.4.** *The value function update rule, $Q_{n+1} = Q_n + \alpha_n(g(Q_n) + M_{n+1})$, has a globally*
891 *asymptotically stable equilibrium, $q_*$.*

892 *Proof.* This was shown in Theorem B.2 of Wan et al. (2021).
893 $\square$

894 **Lemma B.5.** *The subtask update rules (B.58) each have a globally asymptotically stable equilib-*
895 *rium, $z_{i_*}$.*

896 *Proof.* Applying the results of Theorems B.1 and B.2 of Wan et al. (2021) under the two-timescales
897 argument, we have that $g(Q_n) \to 0$, $\bar{R}_n \to \bar{r}_*$, that $\{M_n\}$ is a martingale difference sequence, such
898 that $\mathbb{E}[M_{n+1} \mid \mathcal{F}_n] = 0$ a.s., $n \geq 0$, and that $\{M_n\}$ is square-integrable, such that $\mathbb{E}[||M_{n+1}||^2 \mid$
899 $\mathcal{F}_n] \leq K(1 + ||Q_n||^2)$ a.s., $n \geq 0$, for some constant $K > 0$. Given these results, the remaining
900 $\tilde{R}_n(s,a) - \bar{R}_n = \tilde{R}_n(s,a) - \bar{r}_*$ term in the subtask updates (B.58) can be interpreted as a martingale
901 difference sequence, $\{M_n^r\}$, such that $\mathbb{E}[M_{n+1}^r \mid \mathcal{F}_n] = \mathbb{E}[\tilde{R}_{n+1}(s,a) - \bar{r}_* \mid \mathcal{F}_n] = \mathbb{E}[\tilde{R}_{n+1}(s,a) \mid$
902 $\mathcal{F}_n] - \bar{r}_* = 0$ a.s., $n \geq 0$. As such, given Assumptions B.4, B.8, and B.11, to show that the
903 subtask update rules (B.58) each have a globally asymptotically stable equilibrium, we only need to
904 show that the martingale difference sequence, $\{M_n^r\}$, is square-integrable, such that $\mathbb{E}[(M_{n+1}^r)^2 \mid$
905 $\mathcal{F}_n] < \infty$ a.s., $n \geq 0$. Indeed, because $\tilde{R}_n(s,a)$ is bounded, it directly follows that its variance,
906 $\mathbb{E}[(\tilde{R}_n(s,a) - \bar{r}_*)^2]$, is bounded, and as such, we have that the martingale difference sequence,
907 $\{M_n^r\}$, is square-integrable. Hence, we can conclude that the subtask update rules (B.58) each have
908 a globally asymptotically stable equilibrium, $z_{i_*}$.
909 $\square$

**Lemma B.6.** $\sup_n(||Q_n|| + ||Z_n||) < \infty$ *a.s.*

*Proof.* It was shown in Theorem B.2 of Wan et al. (2021) that $\sup_n(||Q_n||) < \infty$ a.s. Hence, we only need to show that $\sup_n(||Z_n||) < \infty$ a.s. To this end, we can apply Theorem 7 in Section 3 of Borkar (2009). This theorem requires 4 assumptions:

- **(A1)** The function $g$ is Lipschitz: $||g(x) - g(y)|| \leq L||x - y||$ for some $0 < L < \infty$.

- **(A2)** The sequence $\{\eta_{z_i}\alpha_n\}$ satisfies $\eta_{z_i}\alpha_n > 0$, $\sum \eta_{z_i}\alpha_n = \infty$, and $\sum \eta_{z_i}^2\alpha_n^2 < \infty$.

- **(A3)** $\{M_n\}$ and $\{M_n^r\}$ are martingale difference sequences that are square-integrable.

- **(A4)** The functions $g_d(x) \doteq g(dx)/d$, $d \geq 1, x \in \mathbb{R}^k$, satisfy $g_d(x) \to g_*(x)$ as $d \to \infty$, uniformly on compacts for some $g_* \in C(\mathbb{R}^k)$. Furthermore, the ODE $\dot{x}_t = g_*(x_t)$ has the origin as its unique globally asymptotically stable equilibrium.

Under the two-timescales argument, the results of Theorems B.1 and B.2 of Wan et al. (2021) apply, thereby satisfying the above assumptions, except for the assumptions regarding $\{\eta_{z_i}\alpha_n\}$ and $\{M_n^r\}$. In this regard, Assumptions B.4 and B.8 satisfy Assumption **(A2)**. Moreover, we showed in Lemma B.5 that $\{M_n^r\}$ is indeed a martingale difference sequence that is square-integrable. As such, Assumptions **(A1)** - **(A4)** are verified, meaning that we can apply the results of Theorem 7 in Section 3 of Borkar (2009) to conclude that $\sup_n(||Z_n||) < \infty$ a.s., and hence, that $\sup_n(||Q_n|| + ||Z_n||) < \infty$ a.s. $\qquad\square$

As such, we have now verified the 3 assumptions required by Theorem 2 in Section 6 of Borkar (2009), which means that we can apply the result of the theorem to conclude that $Z_{i,n} \to z_{i_*} \forall z_i \in \mathcal{Z}$ a.s. as $n \to \infty$.

Finally, as was done in the proof for linear subtask functions, we conclude the proof by considering $z_{i,\pi_t} \forall z_i \in \mathcal{Z}$, where $\pi_t$ is a greedy policy with respect to $Q_t$. Given that $Q_t \to q_*$ and $\bar{r}_{\pi_t} \to \bar{r}_*$ a.s., it directly follows from Definition 5.1 that $z_{i,\pi_t} \to z_{i_*} \forall z_i \in \mathcal{Z}$ a.s.

This completes the proof of Theorem B.2.1.

## C  Leveraging the RED RL Framework for CVaR Optimization

This appendix contains details regarding the adaptation of the RED RL framework for CVaR optimization. We first derive an appropriate subtask function, then use it to adapt the RED RL algorithms (see Appendix A) for CVaR optimization. In doing so, we arrive at the *RED CVaR algorithms*, which are presented in full at the end of this appendix. These RED CVaR algorithms allow us to optimize CVaR (and VaR) without the use of an augmented state-space or an explicit bi-level optimization. We also provide a convergence proof for the tabular RED CVaR Q-learning algorithm, which shows that the VaR and CVaR estimates converge to the optimal long-run VaR and CVaR, respectively.

### C.1  A Subtask-Driven Approach for CVaR Optimization

In this section, we use the RED RL framework to derive a subtask-driven approach for CVaR optimization that does not require an augmented state-space or an explicit bi-level optimization. To begin, let us consider Equation (7), which is displayed below as Equation (C.1) for convenience:

$$\text{CVaR}_\tau(R_t) = \sup_{y \in \mathcal{R}} \mathbb{E}[y - \frac{1}{\tau}(y - R_t)^+] \tag{C.1a}$$

$$= \mathbb{E}[\text{VaR}_\tau(R_t) - \frac{1}{\tau}(\text{VaR}_\tau(R_t) - R_t)^+], \tag{C.1b}$$

where $\tau \in (0, 1)$ denotes the CVaR parameter, and $R_t$ denotes the observed per-step reward.

We can see from Equation (C.1) that CVaR can be interpreted as an expectation (or average) of sorts, which suggests that it may be possible to leverage the average-reward MDP to optimize this expectation, by treating the reward CVaR as the average-reward, $\bar{r}_\pi$, that we want to optimize. However, this requires that we know the optimal value of the scalar, $y$, because the expectation in Equation (C.1b) only holds for this optimal value (which corresponds to the per-step reward VaR). Unfortunately, this optimal value is typically not known beforehand, so in order to optimize CVaR, we also need to optimize $y$.

Importantly, we can utilize RED RL framework to turn the optimization of $y$ into a subtask, such that CVaR is the primary control objective (i.e., the $\bar{r}_\pi$ that we want to optimize), and VaR ($y$ in Equation (C.1)), is the (single) subtask. This is in contrast to existing MDP-based methods, which typically leverage Equation (C.1) when optimizing for CVaR by augmenting the state-space with a state that corresponds (either directly or indirectly) to an estimate of $\text{VaR}_\tau(R_t)$ (in this case, $y$), and solving the bi-level optimization shown in Equation (8), thereby increasing computational costs.

To utilize the RED RL framework, we first need to derive a valid subtask function for CVaR that satisfies the requirements of Definition 5.1. Let us consider Equation (C.1). We can see that if we treat the expression inside the expectation in Equation (C.1) as our subtask function, $f$ (see Definition 5.1), then we have a piecewise linear subtask function that is invertible with respect to each input given all other inputs, where the subtask, VaR, is independent of the observed per-step reward. Hence, we can adapt Equation (C.1) as our subtask function (given that is satisfies Definition 5.1), as follows:

$$\tilde{R}_t = \text{VaR} - \frac{1}{\tau}(\text{VaR} - R_t)^+, \tag{C.2}$$

where $R_t$ is the observed per-step reward, $\tilde{R}_t$ is the extended per-step reward, VaR is the value-at-risk of the observed per-step reward, and $\tau$ is the CVaR parameter. Importantly, this is a valid subtask function with the following properties: the average (or expected value) of the *extended* reward corresponds to the CVaR of the *observed* reward, and the optimal average of the *extended* reward corresponds to the optimal CVaR of the *observed* reward. This is formalized as Corollaries C.1 - C.4 below:

976    **Corollary C.1.** *The function presented in Equation* (C.2) *is a valid subtask function.*

977    *Proof.* The function presented in Equation (C.2) is clearly a piecewise linear function that is invert-
978    ible with respect to each input given all other inputs. Moreover, the subtask, VaR, is independent of
979    the observed per-step reward. Hence, this function satisfies Definition 5.1 for the subtask, VaR.   □

980    **Corollary C.2.** *If the subtask, VaR (from Equation* (C.2)*) is estimated, and such an estimate is equal*
981    *to the long-run VaR of the observed reward, then the average (or expected value) of the extended*
982    *reward, $\tilde{R}_t$, from Equation* (C.2) *is equal to the long-run CVaR of the observed reward.*

983    *Proof.* This follows directly from Equation (C.1b).   □

984    **Corollary C.3.** *If the subtask, VaR (from Equation* (C.2)*) is estimated, and the resulting average*
985    *of the extended reward from Equation* (C.2) *is equal to the long-run CVaR of the observed reward,*
986    *then the VaR estimate is equal to the long-run VaR of the observed reward.*

987    *Proof.* This follows directly from Equation (C.1b).   □

988    **Corollary C.4.** *A policy that yields an optimal long-run average of the extended reward, $\tilde{R}_t$, from*
989    *Equation* (C.2) *is a CVaR-optimal policy. In other words, the optimal long-run average of the*
990    *extended reward corresponds to the optimal long-run CVaR of the observed reward.*

991    *Proof.* For a given policy, we know from Equation (C.1a) that, across a range of VaR estimates, the
992    best possible long-run average of the extended reward for that policy corresponds to the long-run
993    CVaR of the observed reward for that same policy. Hence, the best possible long-run average of the
994    extended reward that can be achieved across various policies and VaR estimates, corresponds to the
995    optimal long-run CVaR of the observed reward.   □

996    As such, we now have a valid subtask function with a subtask, VaR, and an extended reward whose
997    average, when optimized, corresponds to the optimal CVaR of the observed reward. We are now
998    ready to apply the RED RL framework. First, we can derive the reward-extended TD error update
999    for our subtask, VaR, using the methodology outlined in Section 5.1, where, in this case, we have a
1000   piecewise linear subtask function with two segments. The resulting subtask update is as follows:

$$\text{VaR}_{t+1} = \begin{cases} \text{VaR}_t + \eta\alpha_t\left(\delta_t + \text{CVaR}_t - \text{VaR}_t\right), & R_{t+1} \geq \text{VaR}_t \\ \text{VaR}_t + \eta\alpha_t\left(\left(\frac{\tau}{\tau-1}\right)\delta_t + \text{CVaR}_t - \text{VaR}_t\right), & R_{t+1} < \text{VaR}_t \end{cases}, \tag{C.3}$$

1001   where $\delta_t$ is the regular TD error, and $\eta\alpha_t$ is the step size.

1002   With this update, we now have all the components needed to utilize the RED algorithms in Appendix
1003   A to optimize CVaR (where CVaR corresponds to the $\bar{r}_\pi$ that we want to optimize). We call these
1004   CVaR-specific algorithms, the *RED CVaR algorithms*. The full algorithms are included at the end of
1005   this appendix.

1006   We now present the tabular *RED CVaR Q-learning* algorithm, along with a convergence proof which
1007   shows that the VaR and CVaR estimates converge to the optimal long-run VaR and CVaR of the
1008   observed reward, respectively:

**RED CVaR Q-learning algorithm (tabular):** We update a table of estimates, $Q_t : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ as follows:

$$\tilde{R}_{t+1} = \text{VaR}_t - \frac{1}{\tau}(\text{VaR}_t - R_{t+1})^+ \tag{C.4a}$$

$$\delta_t = \tilde{R}_{t+1} - \text{CVaR}_t + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t) \tag{C.4b}$$

$$Q_{t+1}(S_t, A_t) = Q_t(S_t, A_t) + \alpha_t \delta_t \tag{C.4c}$$

$$Q_{t+1}(s, a) = Q_t(s, a), \quad \forall s, a \neq S_t, A_t \tag{C.4d}$$

$$\text{CVaR}_{t+1} = \text{CVaR}_t + \eta_{\text{CVaR}} \alpha_t \delta_t \tag{C.4e}$$

$$\text{VaR}_{t+1} = \begin{cases} \text{VaR}_t + \eta_{\text{VaR}} \alpha_t \left(\delta_t + \text{CVaR}_t - \text{VaR}_t\right), & R_{t+1} \geq \text{VaR}_t \\ \text{VaR}_t + \eta_{\text{VaR}} \alpha_t \left(\left(\frac{\tau}{\tau-1}\right)\delta_t + \text{CVaR}_t - \text{VaR}_t\right), & R_{t+1} < \text{VaR}_t \end{cases}, \tag{C.4f}$$

where $R_t$ is the observed reward, $\text{VaR}_t$ is the VaR estimate, $\text{CVaR}_t$ is the CVaR estimate, $\alpha_t$ is the step size, $\delta_t$ is the TD error, and $\eta_{\text{CVaR}}, \eta_{\text{VaR}}$ are positive scalars.

**Theorem C.1.1.** *The RED CVaR Q-learning algorithm* (C.4) *converges, almost surely, $\text{CVaR}_t$ to $\text{CVaR}^*$, $\text{VaR}_t$ to $\text{VaR}^*$, $\text{CVaR}_{\pi_t}$ to $\text{CVaR}^*$, $\text{VaR}_{\pi_t}$ to $\text{VaR}^*$, and $Q_t$ to a solution of $q$ in the Bellman Equation* (4), *up to an additive constant, $c$, where $\pi_t$ is any greedy policy with respect to $Q_t$, if the following assumptions hold: 1) the MDP is communicating, 2) the solution of $q$ in* (4) *is unique up to a constant, 3) the step sizes are decreased appropriately as per Assumptions B.3 and B.4, 4) all the state–action pairs are updated an infinite number of times, 5) the ratio of the update frequency of the most-updated state–action pair to the least-updated state–action pair is finite, 6) the subtask function outlined in Equation* (C.2) *is in accordance with Definition 5.1, and 7) $\eta_{\text{VaR}} \alpha_t$ decreases to 0 appropriately, as per Assumption B.8.*

*Proof.* By definition, the RED CVaR Q-learning algorithm (C.4) is of the form of the generic RED Q-learning algorithm (16), where $\text{CVaR}_t$ corresponds to $\bar{R}_t$ and $\text{VaR}_t$ corresponds to $Z_{i,t}$ for a single subtask. We also know from Corollary C.1 that the subtask function used is valid. Hence, Theorem 5.3 applies, such that:

*i)* $\text{CVaR}_t$ and $\text{CVaR}_{\pi_t}$ converge a.s. to the optimal long-run average, $\bar{r}*$, of the extended reward from the subtask function (i.e., the optimal long-run average of $\tilde{R}_t$),

*ii)* $\text{VaR}_t$ and $\text{VaR}_{\pi_t}$ converge a.s. to the corresponding optimal subtask value, $z*$, and

*iii)* $Q_t$ converges to a solution of $q$ in the Bellman Equation (4),

all up to an additive constant, $c$.

Hence, to complete the proof, we need to show that $\bar{r}* = \text{CVaR}^*$ and $z* = \text{VaR}^*$:

From Corollary C.4 we know that the optimal long-run average of the extended reward corresponds to the optimal long-run CVaR of the observed reward, hence we can conclude that $\bar{r}* = \text{CVaR}^*$. Finally, from Corollary C.3 we can deduce that since $\text{CVaR}_t$ converges a.s. to $\text{CVaR}^*$, then $z*$ must correspond to $\text{VaR}^*$.

This completes the proof. $\qquad\square$

As such, with the RED CVaR Q-learning algorithm, we now have a way to optimize the long-run CVaR (and VaR) of the observed reward without the use of an augmented state-space, or an explicit bi-level optimization. See Section 6 and Appendix D for empirical results obtained when using the RED CVaR algorithms.

## C.2  Additional Commentary

We now provide additional commentary on the subtask-driven approach for CVaR optimization:

**Remark C.1.** *A natural question to ask would be whether we can extend these convergence results to the prediction case. In other words, can we show that a tabular RED CVaR TD-learning algorithm will converge to the long-run VaR and CVaR of the observed reward induced by following a given policy? It turns out that, because we are not optimizing the expectation in Equation* (C.1a) *when doing prediction (we are only learning it), we cannot guarantee that we will eventually find the optimal VaR estimate, which implies that we may not recover the CVaR value (since Equation* (C.1b) *only holds to the optimal VaR value). However, this is not to say that a RED CVaR TD-learning algorithm has no use. In fact, we do use such an algorithm as part of an actor-critic architecture for optimizing CVaR in the inverted pendulum experiment (see Appendix D). Empirically, as discussed in Section 6, we find that this actor-critic approach is able to find the optimal CVaR policy.*

**Remark C.2.** *It should be noted that in the risk measure literature, risk measures are typically classified into two categories:* static *or* dynamic. *This classification is based on the* time consistency *of the risk measure that one aims to optimize* Boda and Filar *(2006). Curiously, in our case the CVaR that we aim to optimize does not fit into either category perfectly. One could make the argument that the CVaR that we aim to optimize most closely matches the* static *category, given that there is some time inconsistency before $t \to \infty$. Conversely, one could make a different argument that the CVaR that we aim to optimize most closely resembles the* dynamic *category since the sum over $t$ for the average-reward is outside of the CVaR operator (see Theorem 1 of* Xia et al. *(2023)), such that an optimal deterministic stationary policy exists (unlike the static case; see* Bäuerle and Ott *(2011)). This does not affect the significance of our results, but rather suggests that a third category of risk measures may be needed to capture such nuances that occur in the average-reward setting.*

1067 **C.3 RED CVaR Algorithms**

1068 Below is the pseudocode for the RED CVaR algorithms.

---

**Algorithm 5** RED CVaR Q-Learning (Tabular)

---

**Input:** the policy $\pi$ to be used (e.g., $\varepsilon$-greedy)
**Algorithm parameters:** step size parameters $\alpha, \eta_{\text{CVaR}}, \eta_{\text{VaR}}$, CVaR parameter $\tau$
Initialize $Q(s, a) \,\forall s, a$ (e.g. to zero)
Initialize CVaR arbitrarily (e.g. to zero)
Initialize VaR arbitrarily (e.g. to zero)
Obtain initial $S$
**while** still time to train **do**
   $A \leftarrow$ action given by $\pi$ for $S$
   Take action $A$, observe $R, S'$
   $\tilde{R} = \text{VaR} - \frac{1}{\tau} \max\{\text{VaR} - R, 0\}$
   $\delta = \tilde{R} - \text{CVaR} + \max_a Q(S', a) - Q(S, A)$
   $Q(S, A) = Q(S, A) + \alpha\delta$
   $\text{CVaR} = \text{CVaR} + \eta_{\text{CVaR}}\alpha\delta$
   **if** $R \geq \text{VaR}$ **then**
      $\text{VaR} = \text{VaR} + \eta_{\text{VaR}}\alpha(\delta + \text{CVaR} - \text{VaR})$
   **else**
      $\text{VaR} = \text{VaR} + \eta_{\text{VaR}}\alpha\left(\left(\frac{\tau}{\tau-1}\right)\delta + \text{CVaR} - \text{VaR}\right)$
   **end if**
   $S = S'$
**end while**
return $Q$

---

**Algorithm 6** RED CVaR Actor-Critic

---

**Input:** a differentiable state-value function parameterization $\hat{v}(s, \boldsymbol{w})$; a differentiable policy parameterization $\pi(a \mid s, \boldsymbol{\theta})$
**Algorithm parameters:** step size parameters $\alpha, \eta_\pi, \eta_{\text{CVaR}}, \eta_{\text{VaR}}$, CVaR parameter $\tau$
Initialize state-value weights $\boldsymbol{w} \in \mathbb{R}^d$ and policy weights $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ (e.g. to $\boldsymbol{0}$)
Initialize CVaR arbitrarily (e.g. to zero)
Initialize VaR arbitrarily (e.g. to zero)
Obtain initial $S$
**while** still time to train **do**
   $A \sim \pi(\cdot \mid S, \boldsymbol{\theta})$
   Take action $A$, observe $R, S'$
   $\tilde{R} = \text{VaR} - \frac{1}{\tau} \max\{\text{VaR} - R, 0\}$
   $\delta = \tilde{R} - \text{CVaR} + \hat{v}(S', \boldsymbol{w}) - \hat{v}(S, \boldsymbol{w})$
   $\boldsymbol{w} = \boldsymbol{w} + \alpha\delta\nabla\hat{v}(S, \boldsymbol{w})$
   $\boldsymbol{\theta} = \boldsymbol{\theta} + \eta_\pi\alpha\delta\nabla\ln\pi(A \mid S, \boldsymbol{\theta})$
   $\text{CVaR} = \text{CVaR} + \eta_{\text{CVaR}}\alpha\delta$
   **if** $R \geq \text{VaR}$ **then**
      $\text{VaR} = \text{VaR} + \eta_{\text{VaR}}\alpha(\delta + \text{CVaR} - \text{VaR})$
   **else**
      $\text{VaR} = \text{VaR} + \eta_{\text{VaR}}\alpha\left(\left(\frac{\tau}{\tau-1}\right)\delta + \text{CVaR} - \text{VaR}\right)$
   **end if**
   $S = S'$
**end while**
return $\boldsymbol{w}, \boldsymbol{\theta}$

---

## D   Numerical Experiments

This appendix contains details regarding the numerical experiments performed as part of this work. We discuss the experiments performed in the *red-pill blue-pill* environment (see Appendix E for more details on the red-pill blue-pill environment), as well as the experiments performed in the *inverted pendulum* environment.
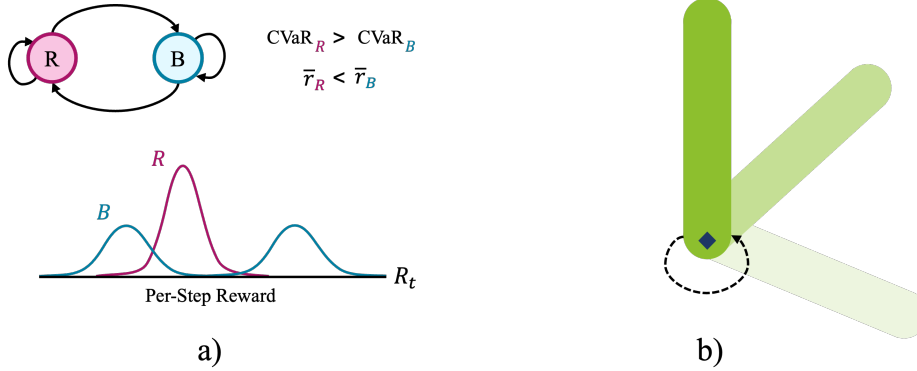


Figure D.1: An illustration of the **a)** red-pill blue-pill, and **b)** inverted pendulum environments.

The aim of the experiments was to contrast and compare the RED RL algorithms (see Appendix C) with the Differential learning algorithms from Wan et al. (2021) in the context of CVaR optimization. In particular, we aimed to show how the RED RL algorithms could be utilized to optimize for CVaR (without the use of an augmented state-space or an explicit bi-level optimization scheme), and contrast the results to those of the Differential learning algorithms, which served as a sort of 'baseline' to illustrate how our *risk-aware* approach contrasts a *risk-neutral* approach. In other words, we aimed to show whether our algorithms could successfully enable a learning agent to act in a risk-aware manner instead of the usual risk-neutral manner.

In terms of the algorithms used, Algorithm 5 corresponds to the RED CVaR Q-learning algorithm used in the red-pill blue-pill experiment, and Algorithm 6 corresponds to the RED CVaR Actor-Critic algorithm used in the inverted pendulum experiment. In terms of the Differential learning algorithms used for comparison (see Appendix D.3 for the full algorithms), Algorithm 7 corresponds to the Differential Q-learning algorithm used in the red-pill blue-pill experiment, and Algorithm 8 corresponds to the Differential Actor-Critic algorithm used in the inverted pendulum experiment.

### D.1   Red-Pill Blue-Pill Experiment

In the first experiment, we consider a two-state environment that we created for the purposes of testing our algorithms. It is called the *red-pill blue-pill* environment (see Appendix D), where at every time step an agent can take either a 'red pill', which takes them to the 'red world' state, or a 'blue pill', which takes them to the 'blue world' state. Each state has its own characteristic per-step reward distribution, and in this case, for a sufficiently low CVaR parameter, $\tau$, the red world state has a per-step reward distribution with a lower (worse) mean but higher (better) CVaR compared to the blue world state. As such, this task allows us to answer the following question: *can the RED CVaR algorithms successfully get the agent to learn a policy that prioritizes optimizing the reward CVaR over the average-reward?* In particular, we would expect that the RED CVaR algorithms learn a policy that prefers to stay in the red world, and that the (risk-neutral) Differential algorithms (from Wan et al. (2021)) learn a policy that prefers to stay in the blue world. This task is illustrated in Figure D.1a).

1101 For this experiment, we ran both algorithms using various combinations of step sizes for each algo-
1102 rithm. We used an $\varepsilon$-greedy policy with a fixed epsilon of 0.1, and a CVaR parameter, $\tau$, of 0.25.
1103 We set all initial guesses to zero. We ran the algorithms for 100k time steps.

1104 For the Differential Q-learning algorithm, we tested every combination of the value function step
1105 size, $\alpha \in \{2e\text{-}1, 2e\text{-}2, 2e\text{-}3, 2e\text{-}4, 1/n\}$ (where $1/n$ refers to a step size sequence that decreases
1106 the step size according to the time step, $n$), with the average-reward step size, $\eta\alpha$, where $\eta \in$
1107 $\{1e\text{-}4, 1e\text{-}3, 1e\text{-}2, 1e\text{-}1, 1.0, 2.0\}$, for a total of 30 unique combinations. Each combination was run
1108 50 times using different random seeds, and the results were averaged across the runs. The resulting
1109 (averaged) average-reward over the last 1,000 time steps is displayed in Figure D.2. As shown in
1110 the figure, a value function step size of 2e-4 and an average-reward $\eta$ of 1.0 resulted in the highest
1111 average-reward in the final 1,000 time steps in the red-pill blue-pill task. These are the parameters
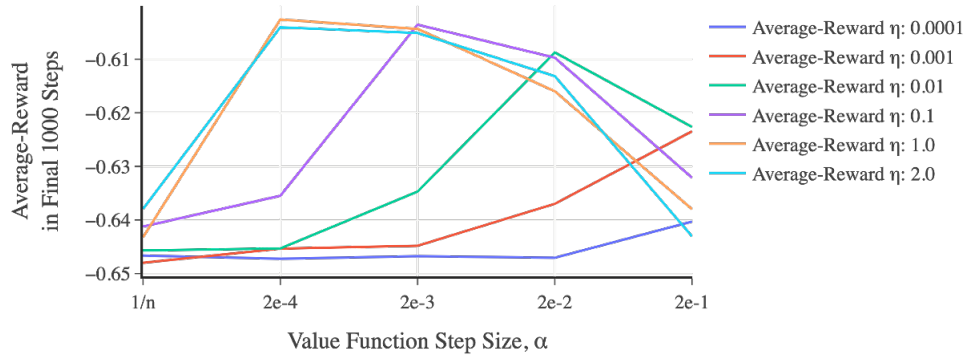1112 used to generate the results displayed in Figure 3a).



Figure D.2: Step size tuning results for the red-pill blue-pill task when using the Differential Q-learning algorithm. The average-reward in the final 1,000 steps is displayed for various combinations of value function and average-reward step sizes.

1113 For the RED CVaR Q-learning algorithm, we tested every combination of the value function
1114 step size, $\alpha \in \{2e\text{-}1, 2e\text{-}2, 2e\text{-}3, 2e\text{-}4, 1/n\}$, with the average-reward (in this case CVaR)
1115 $\eta \in \{1e\text{-}4, 1e\text{-}3, 1e\text{-}2, 1e\text{-}1, 1.0, 2.0\}$, and the VaR $\eta \in \{1e\text{-}4, 1e\text{-}3, 1e\text{-}2, 1e\text{-}1, 1.0, 2.0\}$, for a total
1116 of 180 unique combinations. Each combination was run 50 times using different random seeds, and
1117 the results were averaged across the runs. A value function step size of 2e-2, an average-reward
1118 (CVaR) $\eta$ of 1e-1, and a VaR $\eta$ of 1e-1 yielded the best results and were used to generate the results
1119 displayed in Figures 3a) and 4a).
1120

**Follow-up Experiment: Varying the CVaR Parameter**

1122 Given the results shown in Figure 3a), we can see that, with proper hyperparameter tuning, the
1123 tabular RED CVaR Q-learning algorithm is able to reliably find the optimal CVaR policy for a CVaR
1124 parameter, $\tau$, of 0.25. In the context of the red-pill blue-pill environment, this means that the agent
1125 learns to stay in the red world state because the state has a characteristic reward distribution with
1126 a better (higher) CVaR compared to the blue world state. By contrast, the risk-neutral differential
1127 algorithm yields an average-reward optimal policy that dictates that the agent should stay in the blue
1128 world state because the state has a better (higher) average reward compared to the red world state.

1129 Now consider what would happen if we used the RED CVaR Q-learning algorithm with a $\tau$ of 0.99.
1130 By definition, a CVaR corresponding to a $\tau \approx 1.0$ is equivalent to the average reward. Hence, with
1131 a $\tau$ of 0.99, we would expect that the optimal CVaR policy corresponds to staying in the blue world
1132 state (since it has the better average reward). This means that for some $\tau$ between 0.25 and 0.99,
1133 there is a critical point where the CVaR-optimal policy changes from staying in the red world (let us
1134 call this the *red policy*) to staying in the blue world state (let us call this the *blue policy*).

1135  We can estimate this critical point using simple Monte Carlo (MC). We are able to use MC in this
1136  case because both policies effectively stay in a single state (the red or blue world state), such that
1137  the CVaR of the policies can be estimated by sampling the characteristic reward distribution of each
1138  state, while accounting for the exploration $\varepsilon$. Figure D.3 shows the MC estimate of the CVaR of the
1139  red and blue policies for a range of CVaR parameters, assuming an exploration $\varepsilon$ of 0.1. Note that
1140  we used the same distribution parameters listed in Appendix E for the red-pill blue-pill environment.
1141  As shown in Figure D.3, this critical point occurs somewhere around $\tau \approx 0.8$.
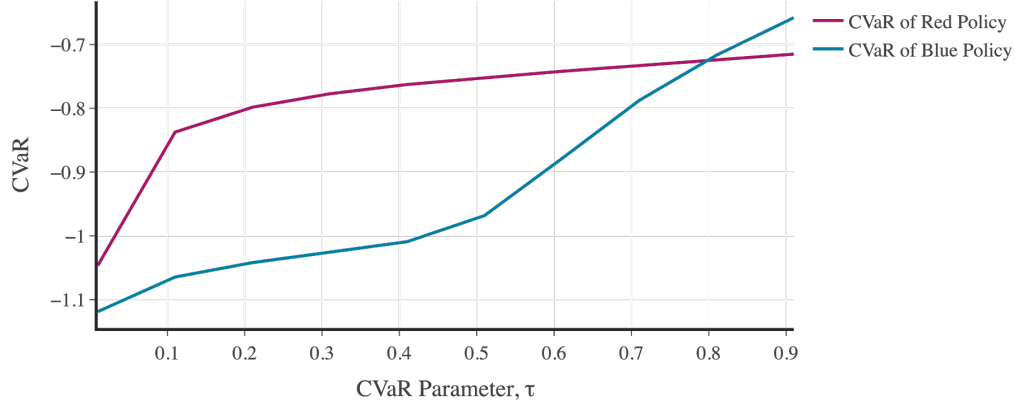


Figure D.3: Monte Carlo estimates of the CVaR of the red and blue policies for a range of CVaR parameters in the red-pill blue-pill environment.

1142  Hence, one way that we can further validate the tabular RED CVaR Q-learning algorithm, is by
1143  re-running the red-pill blue-pill experiment for different CVaR parameters, and seeing if the optimal
1144  CVaR policy indeed changes at a $\tau \approx 0.8$. Importantly, this allows us to empirically validate
1145  whether the algorithm actually optimizes at the desired risk level. When running this experiment,
1146  we used the same hyperparameters used to generate the results in Figure 3a). We ran the experiment
1147  for $\tau \in \{0.1, 0.25, 0.5, 0.75, 0.85, 0.9\}$. For each $\tau$, we performed 50 runs using different random
1148  seeds, and the results were averaged across the runs.
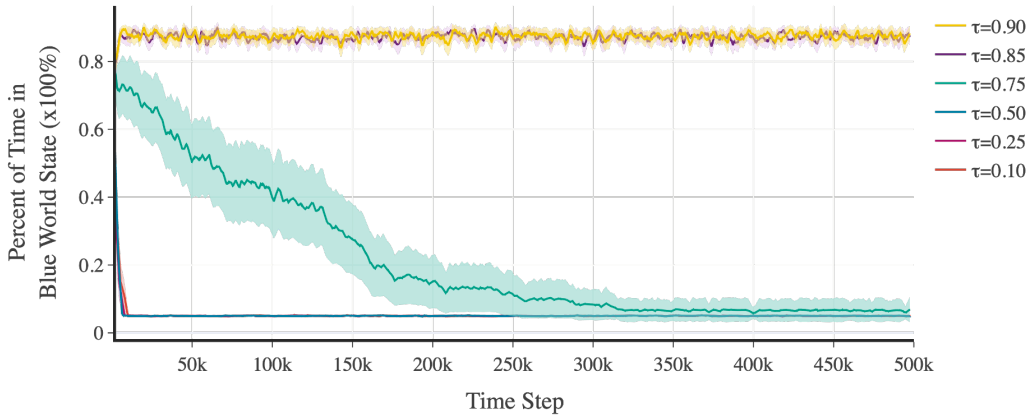


Figure D.4: Rolling percent of time that the agent stays in the blue world state as learning progresses when using the RED CVaR Q-learning algorithm in the red-pill blue-pill environment for a range of CVaR parameters. A solid line denotes the mean percent of time spent in the blue world state, and the corresponding shaded region denotes a 95% confidence interval over 50 runs.

1149 Figure D.4 shows the results of this experiment. In particular, the figure shows a rolling percent
1150 of time that the agent stays in the blue world state as learning progresses (note that we used an
1151 exploration $\varepsilon$ of 0.1). From the figure, we can see that for $\tau \in \{0.1, 0.25, 0.5, 0.75\}$, the agent
1152 learns to stay in the red world state, and for $\tau \in \{0.85, 0.9\}$, the agent learns to stay in the blue
1153 world state. This is consistent with what we would expect, given that the critical point is $\tau \approx 0.8$.
1154 Hence, these results further validate that our algorithm is able to optimize at the desired risk level.

### D.2 Inverted Pendulum Experiment

1156 In the second experiment, we consider the well-known *inverted pendulum* task, where an agent
1157 learns how to optimally balance an inverted pendulum. We chose this task because it provides
1158 us with the opportunity to test our algorithms in an environment where: 1) we must use function
1159 approximation (given the high-dimensional state-space), and 2) where the optimal CVaR policy and
1160 the optimal average-reward policy is the same policy (i.e., the policy that best balances the pendulum
1161 will yield a limiting reward distribution with both the optimal average-reward and reward CVaR).
1162 This hence allows us to directly compare the performance of our RED algorithms to that of the
1163 regular Differential learning algorithms, as well as to gauge how function approximation affects the
1164 performance of our algorithms. For this task, we utilized a simple actor-critic architecture (Barto
1165 et al., 1983; Sutton and Barto, 2018) as this allowed us to compare the performance of a (non-tabular)
1166 RED TD-learning algorithm with a (non-tabular) Differential TD-learning algorithm. This task is
1167 illustrated in Figure D.1b).

1168 For this experiment, we ran both algorithms using various combinations of step sizes for each algo-
1169 rithm. We used a fixed CVaR parameter, $\tau$, of 0.1. We set all initial guesses to zero. We ran the
1170 algorithms for 100k time steps. For simplicity, we used tile coding (Sutton and Barto, 2018) for both
1171 the value function and policy parameterizations, where we parameterized a softmax policy. For each
1172 parameterization, we used 32 tilings, each with 8 X 8 tiles. By using a linear function approximator
1173 (i.e., tile coding), the gradients for the value function and policy parameterizations can be simplified
1174 as follows:

$$\nabla \hat{v}(s, \boldsymbol{w}) = \boldsymbol{x}(s), \tag{D.1}$$

$$\nabla \ln \pi(a \mid s, \boldsymbol{\theta}) = \boldsymbol{x}_h(s, a) - \sum_{\xi \in \mathcal{A}} \pi(\xi \mid s, \boldsymbol{\theta}) \boldsymbol{x}_h(s, \xi), \tag{D.2}$$

1175 where $s \in \mathcal{S}$, $a \in \mathcal{A}$, $\boldsymbol{x}(s)$ is the state feature vector, and $\boldsymbol{x}_h(s, a)$ is the softmax preference vector.

1176 For the Differential Actor-Critic algorithm, we tested every combination of the value function step
1177 size, $\alpha \in \{2e-2, 2e-3, 2e-4, 1/n\}$, with $\eta$'s for the average-reward and policy step sizes, $\eta \alpha$, where
1178 $\eta \in \{1e-3, 1e-2, 1e-1, 1.0, 2.0\}$, for a total of 100 unique combinations. Each combination was
1179 run 10 times using different random seeds, and the results were averaged across the runs. A value
1180 function step size of 2e-3, a policy $\eta$ of 2.0, and an average-reward $\eta$ of 1e-2 yielded the best results
1181 and were used to generate the results displayed in Figure 3b).

1182 For the RED CVaR Actor-Critic algorithm, we tested every combination of the value function step
1183 size, $\alpha \in \{2e-2, 2e-3, 2e-4, 1/n\}$ (where $1/n$ refers to a step size sequence that decreases the step
1184 size according to the time step, $n$), with $\eta$'s for the average-reward, VaR, and policy step sizes, $\eta \alpha$,
1185 where $\eta \in \{1e-3, 1e-2, 1e-1, 1.0, 2.0\}$, for a total of 500 unique combinations. Each combination
1186 was run 10 times using different random seeds, and the results were averaged across the runs. A
1187 value function step size of 2e-3, a policy $\eta$ of 1e-1, an average-reward (CVaR) $\eta$ of 1e-2, and a VaR
1188 $\eta$ of 1e-2 were used to generate the results displayed in Figures 3b) and 4b).

### D.3 Risk-Neutral Differential Algorithms

Below is the pseudocode for the risk-neutral differential algorithms used for comparison in our experiments.

---

**Algorithm 7** Differential Q-Learning (Tabular)

---

**Input:** the policy $\pi$ to be used (e.g., $\varepsilon$-greedy)
**Algorithm parameters:** step size parameters $\alpha, \eta$
Initialize $Q(s, a) \, \forall s, a$ (e.g. to zero)
Initialize $\bar{R}$ arbitrarily (e.g. to zero)
Obtain initial $S$
**while** still time to train **do**
 $A \leftarrow$ action given by $\pi$ for $S$
 Take action $A$, observe $R, S'$
 $\delta = R - \bar{R} + \max_a Q(S', a) - Q(S, A)$
 $Q(S, A) = Q(S, A) + \alpha\delta$
 $\bar{R} = \bar{R} + \eta\alpha\delta$
 $S = S'$
**end while**
return $Q$

---

**Algorithm 8** Differential Actor-Critic

---

**Input:** a differentiable state-value function parameterization $\hat{v}(s, \boldsymbol{w})$; a differentiable policy parameterization $\pi(a \mid s, \boldsymbol{\theta})$
**Algorithm parameters:** step size parameters $\alpha, \eta_\pi, \eta_{\bar{R}}$
Initialize state-value weights $\boldsymbol{w} \in \mathbb{R}^d$ and policy weights $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ (e.g. to $\boldsymbol{0}$)
Initialize $\bar{R}$ arbitrarily (e.g. to zero)
Obtain initial $S$
**while** still time to train **do**
 $A \sim \pi(\cdot \mid S, \boldsymbol{\theta})$
 Take action $A$, observe $R, S'$
 $\delta = R - \bar{R} + \hat{v}(S', \boldsymbol{w}) - \hat{v}(S, \boldsymbol{w})$
 $\boldsymbol{w} = \boldsymbol{w} + \alpha\delta\nabla\hat{v}(S, \boldsymbol{w})$
 $\boldsymbol{\theta} = \boldsymbol{\theta} + \eta_\pi\alpha\delta\nabla\ln\pi(A \mid S, \boldsymbol{\theta})$
 $\bar{R} = \bar{R} + \eta_{\bar{R}}\alpha\delta$
 $S = S'$
**end while**
return $\boldsymbol{w}, \boldsymbol{\theta}$

---

## E   Red-Pill Blue-Pill Environment

This appendix contains a Python implementation of the *red-pill blue-pill* environment introduced in this work. The environment consists of a two-state MDP, where at every time step an agent can take either a 'red pill', which takes them to the 'red world' state, or a 'blue pill', which takes them to the 'blue world' state. Each state has its own characteristic per-step reward distribution, and in this case, for a sufficiently low CVaR parameter, $\tau$, the red world state has a per-step reward distribution with a lower (worse) mean but higher (better) CVaR compared to the blue world state. More specifically, the red world state reward distribution is characterized as a gaussian distribution with a mean of -0.7 and a standard deviation of 0.05. The blue world state is characterized by a mixture of two gaussian distributions with means of -1.0 and -0.2, and standard deviations of 0.05. We assume all rewards are non-positive. The Python implementation of the environment is provided below:

```python
import pandas as pd
import numpy as np

class EnvironmentRedPillBluePill:
  def __init__(self, dist_2_mix_coefficient=0.5):
    # set distribution parameters
    self.dist_1 = {'mean': -0.7, 'stdev': 0.05}
    self.dist_2a = {'mean': -1.0, 'stdev': 0.05}
    self.dist_2b = {'mean': -0.2, 'stdev': 0.05}
    self.dist_2_mix_coefficient = dist_2_mix_coefficient

    # start state
    self.start_state = np.random.choice(['redworld', 'blueworld'])

  def env_start(self, start_state=None):
    # return initial state
    if pd.isnull(start_state):
      return self.start_state
    else:
      return start_state

  def env_step(self, state, action, terminal=False):
    if action == 'red_pill':
      next_state = 'redworld'
    elif action == 'blue_pill':
      next_state = 'blueworld'

    if state == 'redworld':
      reward = np.random.normal(loc=self.dist_1['mean'],
                                scale=self.dist_1['stdev'])
    elif state == 'blueworld':
      dist = np.random.choice(['dist2a', 'dist2b'],
                              p=[self.dist_2_mix_coefficient,
                                 1 - self.dist_2_mix_coefficient])
      if dist == 'dist2a':
        reward = np.random.normal(loc=self.dist_2a['mean'],
                                  scale=self.dist_2a['stdev'])
      elif dist == 'dist2b':
        reward = np.random.normal(loc=self.dist_2b['mean'],
                                  scale=self.dist_2b['stdev'])

    return min(0, reward), next_state, terminal
```