

DUAL-OBJECTIVE REINFORCEMENT LEARNING WITH NOVEL HAMILTON-JACOBI-BELLMAN FORMULATIONS

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
Anonymous authors
Paper under double-blind review

ABSTRACT

Hard constraints in reinforcement learning (RL) often degrade policy performance. Lagrangian methods offer a way to blend objectives with constraints, but require intricate reward engineering and parameter tuning. In this work, we extend recent advances that connect Hamilton-Jacobi (HJ) equations with RL to propose two novel value functions for dual-objective satisfaction. Namely, we address: 1) the **Reach-Always-Avoid** (RAA) problem – of achieving distinct reward and penalty thresholds – and 2) the **Reach-Reach** (RR) problem – of achieving thresholds of two distinct rewards. In contrast with temporal logic approaches, which typically involve representing an automaton, we derive explicit, tractable Bellman forms in this context via decomposition. Specifically, we prove that the RAA and RR problems may be rewritten as compositions of previously studied HJ-RL problems. We leverage our analysis to propose a variation of Proximal Policy Optimization (**DO-HJ-PPO**), and demonstrate that it produces distinct behaviors from previous approaches, out-competing a number of baselines in success, safety and speed across a range of tasks for safe-arrival and multi-target achievement.

1 INTRODUCTION

The development of special Bellman equations from the Hamilton-Jacobi (HJ) perspective of dynamic programming (DP) has illustrated a novel route to safety and target-achievement in reinforcement learning (RL) Fisac et al. (2019); Hsu et al. (2021). In comparison with the canonical RL discounted-sum cost and corresponding additive DP update, these equations, namely the Safety Bellman Equation (SBE) and Reach-Avoid Bellman Equation (RABE), propagate the minimum (worst) penalty and maximum (best) reward, yielding a value function defined by the outlying performance of a trajectory. In mission-critical applications, where avoiding failure is a necessary condition, these equations have proved invaluable in the field of safe control Mitchell et al. (2005); Ames et al. (2016). By focusing on extremal values rather than discounted sums, the HJ-RL equations induce behaviors that act with respect to the best or worst outcomes in time-optimal fashions, performing far more safely than Lagrangian methods Ganai et al. (2023); So et al. (2024). Accordingly, these updates yield policies with significantly improved performance in target-achievement and obstacle-avoidance tasks over long horizons Yu et al. (2022a;b), relevant to fundamental and practical problems in many domains.

In this work, we advance the existing HJ-RL formulations by generalizing them to compositional problems. To date, the HJ-RL Bellman equations are limited to three operations: Reach (R), wherein the agent seeks to reach a goal (achieve a reward threshold), Avoid (A), wherein the agent seeks to avoid an obstacle (avoid a penalty threshold), and Reach-Avoid (RA), where the agent reaches a goal while avoiding obstacles on the way. In this light, we extend the HJ-RL Bellman equations to two complementary problems concerned with dual-satisfaction, namely the **Reach-Reach** (RR) problem for reaching two goals and the **Reach-Always-Avoid** (RAA) problem for continuing to avoid hazards after reaching a goal, demonstrated in Figure 1. We prove that the RAA and RR have a fundamental structure such that their Bellman equations may be decomposed into combinations of SBE’s and RABE’s. From this theory, we devise **DO-HJ-PPO**, a novel algorithm for learning the RAA and RR values which bootstraps concurrently solved decompositions for coupling on-policy PPO roll-outs. Notably, this allows a user to automatically learn to satisfy dual-objective tasks, for example, in the RAA, the F16 learns to fly into the desired airspace without crashing afterward (Figure 1, top middle-left), and in the RR case, the Hopper learns to jump into a target without diving so it may then achieve the second target (Figure 1, bottom left). The RAA and RR problems are distinct from both standard sum-of-reward values and the simpler HJ-RL formulations, providing new perspectives and performant tools for constrained decision-making.

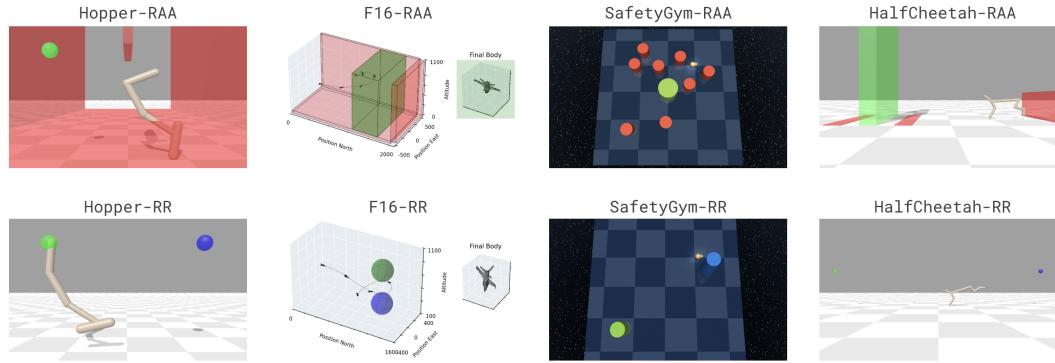


Figure 1: **Depiction of the Reach-Always-Avoid (RAA) and Reach-Reach (RR) Tasks** In the RAA tasks, the zero-level set of the rewards (goals) and penalties (obstacles) are depicted in green and red respectively, while in the RR problem, the zero-level set of the two rewards (two goals) are depicted in green and blue. The RAA value is defined by the minimum of the minimum penalty and maximum reward, inducing the agents to enter the goals at some time without ever entering the obstacles. The RR value is defined by the minimum of the two maximum rewards, inducing the agents to enter both goals at some time.

Our contributions include:

- We introduce novel value functions corresponding to the RAA and RR problems.
- We prove that these value functions and their optimal policies can be decomposed into reach, avoid, and reach-avoid value functions (Theorems 1 and 2).
- We demonstrate the nature of the RAA and RR values and their optimal policies in a simple grid-world example with deep Q -learning (DQN) (Figure 2).
- We propose **DO-HJ-PPO** to solve these value functions, which bootstraps concurrently solved decompositions for effectively coupling the on-policy rollouts (Section 7.2).
- In continuous control tasks, we showcase that with *little to no tuning*, **DO-HJ-PPO** is more successful, safer and faster than Lagrangian and existing HJ-RL baselines (Figure 4).

2 RELATED WORKS

This work involves aspects of safety (e.g. hazard avoidance), liveness (e.g. goal reaching), and balancing competing objectives. We summarize the relevant related works here.

Constrained and Multi-Objective RL. Constrained Markov decision processes (CMDPs) maximize the expected sum of discounted rewards subject to an expected sum of discounted costs, or an instantaneous safety violation function remaining below a set threshold Altman (2021); Achiam et al. (2017a); Wachi and Sui (2020). CMDPs are an effective way to incorporate state constraints into RL problems, and the efficient and accurate solution of the underlying optimization problem has been extensively researched, first by Lagrangian methods and later by an array of more sophisticated techniques Stooke et al. (2020); Li et al. (2024); Chen et al. (2021); Miryoosefi and Jin (2021); Yang et al. (2020). Multi-objective RL is an approach to designing policies that obtain *Pareto-optimal* expected sums of discounted *vector-valued* rewards Wiering et al. (2014); Van and Nowé (2014); Cai et al. (2023), including by deep-Q and other deep learning techniques Mossalam et al. (2016); Abels et al. (2019); Yang et al. (2019). By contrast, this work explicitly balances rewards and penalties in a way that does not require specifying a Lagrange multiplier or similar hyperparameter. Moreover, our work treats goal-reaching and hazard-avoidance as hard constraints, and the learned value function has a direct interpretation in terms of the constraint satisfaction.

Goal-Conditioned RL (GCRL). GCRL simultaneously learns optimal policies for a range of different (but typically related) tasks Liu et al. (2022); Plappert et al. (2018); Ren et al. (2019); Ma et al. (2022); Campero et al. (2020); Trott et al. (2019); Eysenbach et al. (2022); Ma et al. (2022); Campero et al. (2020). In GCRL, states are augmented with information on the current goal. While these goals are in their simplest form mostly independent, some work extends GCRL to more sophisticated composite tasks Chane-Sane et al. (2021). Our work primarily focuses on composing specific learned tasks rather than learning general tasks simultaneously.

108 **Linear Temporal Logic (LTL), Automatic State Augmentation, and Automatons.** Many works
 109 have been explored that merge LTL and RL, canonically focused on Non-Markovian Reward Decision
 110 Processes (NMRDPs) Bacchus et al. (1996). Here, the reward gained at each time step may depend
 111 on the previous state history. Many of these works convert these NMRDPs to MDPs via state
 112 augmentation Bacchus et al. (1997); Thiebaux et al. (2006); Camacho et al. (2021); Icarte et al.
 113 (2018); Camacho et al. (2019). Often the augmented states are taken to be products between an
 114 ordinary state and an automaton state, where the automaton is used to determine "where" in the LTL
 115 specification an agent currently is. Other works using RL for LTL tasks involve MDP verification
 116 Brázil et al. (2014), hybrid systems theory Cohen et al. (2023), GCRL with complex LTL tasks
 117 Qiu et al. (2023), almost-sure objective satisfaction Sadigh et al. (2014), incorporating (un)timed
 118 specifications Hamilton et al. (2022), and using truncated LTL Li et al. (2017). While the problems
 119 we attempt to solve (e.g. reaching multiple goals) can be thought of as specific instantiations of LTL
 120 specifications, our approach to solving these problems is fundamentally different from those in this
 121 line of work. Our state augmentation and subsequent decomposition of the problem are performed
 122 in a specific manner to leverage new HJ-based methods on the subproblems. Through our specific
 123 choice of state augmentation, we still prove that we can achieve an optimal policy in theory (and
 approximately so in practice) despite the non-NMRDP setup.

124 **Hamilton-Jacobi (HJ) Methods.** HJ is a dynamic programming-based framework for solving reach,
 125 avoid, and reach-avoid tasks Mitchell et al. (2005); Fisac et al. (2015). The value functions used in
 126 HJ have the advantage of directly specifying desired behavior, so that a positive value corresponds
 127 to task achievement and a negative value corresponds to task failure. Recent works use RL to find
 128 corresponding optimal policies by leveraging the unconventional Bellman updates associated with
 129 these value functions So et al. (2024); Hsu et al. (2021); Fisac et al. (2019). We build on these works
 130 by extending these advancements to more complex tasks, superficially mirroring the progression from
 131 MDPs to NMRDPs in the LTL-RL literature. Additional works merge HJ and RL, but do not concern
 132 themselves with such composite tasks Ganai et al. (2023); Yu et al. (2022a); Zhu et al. (2024).

134 3 PROBLEM DEFINITION

135 Consider a Markov decision process (MDP) $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, f \rangle$ consisting of finite state and action
 136 spaces \mathcal{S} and \mathcal{A} , and *unknown* discrete dynamics f that define the deterministic transition $s_{t+1} =$
 137 $f(s_t, a_t)$. Let an agent interact with the MDP by selecting an action with policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ to yield
 138 a state trajectory s_t^π , i.e. $s_{t+1}^\pi = f(s_t^\pi, \pi(s_t^\pi))$.

139 In this work, we consider the **Reach-Always-Avoid** (RAA) and **Reach-Reach** (RR) problems, which
 140 both involve the composition of two objectives, which are each specified in terms of the best reward
 141 and worst penalty encountered over time. In the RAA problem, let $r, p : \mathcal{S} \rightarrow \mathbb{R}$ represent a reward
 142 to be maximized and a penalty to be minimized. We will let $q = -p$ for mathematical convenience,
 143 but for conceptual ease we recommend the reader think of trying to minimize the largest-over-time
 144 penalty p rather than maximize the smallest-over-time q . In the RR problem, let $r_1, r_2 : \mathcal{S} \rightarrow \mathbb{R}$ be
 145 two distinct rewards to be maximized. The agent's overall objective is to maximize the *worst-case*
 146 outcome between the best-over-time reward and worst-over-time penalty (in RAA) and the two
 147 best-over-time rewards (in RR), i.e.

$$\begin{aligned} \text{(RAA)} & \left\{ \begin{array}{ll} \text{maximize (w.r.t. } \pi) & \min \left\{ \max_t r(s_t^\pi), \min_t q(s_t^\pi) \right\} \\ \text{s.t.} & s_{t+1}^\pi = f(s_t^\pi, \pi(s_t^\pi)), \\ & s_0^\pi = s, \end{array} \right. \\ \text{(RR)} & \left\{ \begin{array}{ll} \text{maximize (w.r.t. } \pi) & \min \left\{ \max_t r_1(s_t^\pi), \max_t r_2(s_t^\pi) \right\} \\ \text{s.t.} & s_{t+1}^\pi = f(s_t^\pi, \pi(s_t^\pi)), \\ & s_0^\pi = s. \end{array} \right. \end{aligned}$$

148 As the names suggest, these optimization problems are inspired by — but not limited to — tasks
 149 involving goal reaching and hazard avoidance. More specifically, the RAA problem is motivated by a
 150 task in which an agent wishes to both reach a goal \mathcal{G} and perennially avoid a hazard \mathcal{H} (even after it
 151 reaches the goal). The RR problem is motivated by a task in which an agent wishes to reach two goals,
 152 \mathcal{G}_1 and \mathcal{G}_2 , in either order. While these problems are thematically distinct, they are mathematically
 153 complementary (differing by a single max/min operation), and hence we tackle them together.

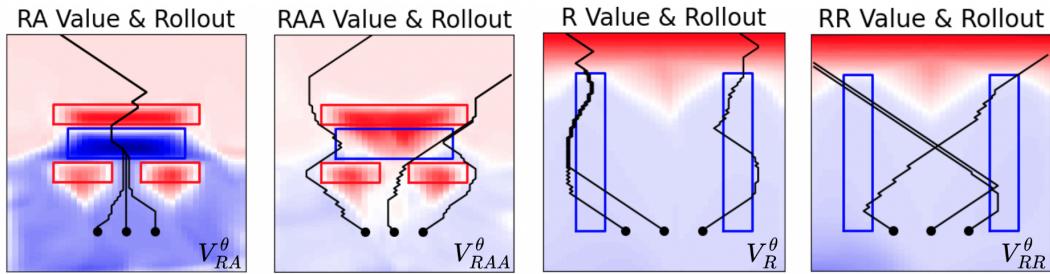


Figure 2: **DQN Grid-World Demonstration of the RAA & RR Problems** We compare our novel formulations with previous HJ-RL formulations (RA & R) in a simple grid-world problem with DQN. The zero-level sets of q (hazards) are highlighted in red, those of r (goals) in blue, and trajectories in black (starting at the dot). In both models, the agents actions are limited to {left, right, straight} and the system flows upwards over time.

The values for any policy in these problems then take the forms V_{RAA}^π and V_{RR}^π ,

$$V_{RAA}^\pi(s) = \min \left\{ \max_t r(s_t^\pi), \min_t q(s_t^\pi) \right\} \quad \text{and} \quad V_{RR}^\pi(s) = \min \left\{ \max_t r_1(s_t^\pi), \max_t r_2(s_t^\pi) \right\}.$$

One may observe that these values are fundamentally different from the infinite-sum value commonly employed in RL Sutton and Barto (2018), and do not accrue over the trajectory but, rather, are determined by certain points. Moreover, while each return considers two objectives, these objectives are combined in worst-case fashion to ensure *dual-satisfaction*. Although many of the works discussed in the previous section approach related tasks (e.g. goal reaching and hazard avoidance) via traditional sum-of-discounted-rewards formulations, these novel value functions have a more direct interpretation in the following sense: if r is positive (only) within \mathcal{G} and q is positive (only) inside \mathcal{H} , $V_{RAA}^\pi(s)$ will be positive if and only if the RAA task will be accomplished by the policy π . Similarly if r_1 and r_2 are positive within \mathcal{G}_1 and \mathcal{G}_2 , respectively, $V_{RR}^\pi(s)$ will be positive if and only if the RR task will be accomplished by the policy π .

4 REACHABILITY AND AVOIDABILITY IN RL

The reach V_R^π , avoid V_A^π , and reach-avoid V_{RA}^π values, respectively defined by

$$V_R^\pi(s) = \max_t r(s_t^\pi), \quad V_A^\pi(s) = \min_t q(s_t^\pi), \quad V_{RA}^\pi(s) = \max_t \min_{\tau \leq t} \left\{ r(s_\tau^\pi), \min_{\tau \leq t} q(s_\tau^\pi) \right\},$$

have been previously studied Fisac et al. (2019) leading to the derivation of special Bellman equations. To put these value functions in context, assume the goal \mathcal{G} is the set of states for which $r(s)$ is positive and the hazard \mathcal{H} is the set of states for which $q(s)$ is non-positive. See Figure 2 for a simple grid-world demonstration comparing the RAA and RR values with the previously existing RA and R values. Then V_R^π , V_A^π , and V_{RA}^π are positive if and only if π causes the agent to eventually reach \mathcal{G} , to always avoid \mathcal{H} , and to reach \mathcal{G} without hitting \mathcal{H} prior to the reach time, respectively. The Reach-Avoid Bellman Equation (RABE), for example, takes the form Hsu et al. (2021)

$$V_{RA}^*(s) = \min \left\{ \max \left\{ \max_{a \in \mathcal{A}} V_{RA}^*(f(s, a)), r(s) \right\}, q(s) \right\},$$

and is associated with optimal policy $\pi_{RA}^*(s)$ (without the need for state augmentation, see the appendix). This formulation does not naturally induce a contraction, but may be discounted to induce contraction by defining $V_{RA}^\gamma(z)$ implicitly via

$$V_{RA}^\gamma(s) = (1 - \gamma) \min\{r(s), q(s)\} + \gamma \min \left\{ \max \left\{ \max_{a \in \mathcal{A}} V_{RA}^\gamma(f(s, a)), r(s) \right\}, q(s) \right\},$$

for each $\gamma \in [0, 1]$. A fundamental result (Proposition 3 in Hsu et al. (2021)) is that

$$\lim_{\gamma \rightarrow 1} V_{RA}^\gamma(s) = V_{RA}(s).$$

These prior value functions and corresponding Bellman equations have proven powerful for these simple reach/avoid/reach-avoid problem formulations. In this work, we generalize the aforementioned results to the broader class involving V_{RAA} (assure no penalty after the reward threshold is achieved) and V_{RR} (achieve multiple rewards optimally). Through this generalization, we are able to train an agent to accomplish more complex tasks with noteworthy performance.

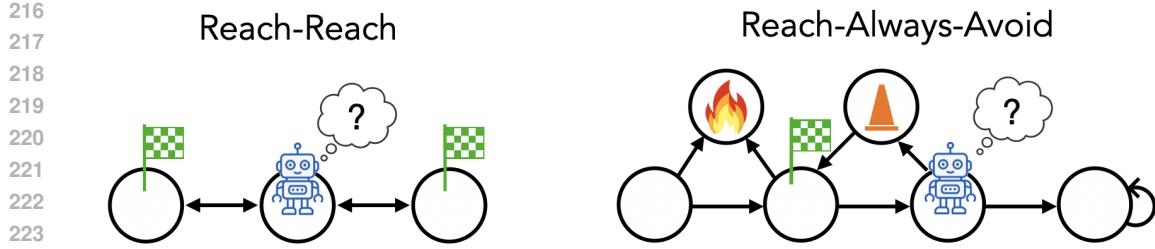


Figure 3: **Examples where a Non-Augmented Policy is Flawed** In both MDPs, consider an agent with no memory. (Left) For a deterministic policy based on the current state, the agent can only achieve one target (RR), as this policy must associate the middle state with either of the two possible actions. (Right) The RAA case is slightly more complex. Assume the robot will make sure to avoid the fire at all costs (which is easily done from the current state). It would also prefer to not encounter the cone hazard, but will do so if needed to achieve the target. From its current state the robot cannot determine whether to pursue the target by crossing the cone or move to the right. The correct decision depends on state history, specifically on whether the robot has already reached the target state or not (e.g. imagine the initial state is on the target state).

5 THE NEED FOR AUGMENTING STATES WITH HISTORICAL INFORMATION

We here discuss a small but important detail regarding the problem formulation. The value functions we introduce may appear similar to the simpler HJ-RL value functions discussed in the previous section; however, in these new formulations the goal of choosing a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ is inherently flawed without state augmentation. In considering multiple objectives over an infinite horizon, situations arise in which the optimal action depends on more than the current state, but rather the **history** the trajectory. This complication is not unique to our problem formulation, but also occurs for NMDPs (see the Related Works section). To those unfamiliar with NMDPs, this at first may seem like a paradox as the MDP is by definition Markov, but the problem occurs not due to the state-transition dynamics but the nature of the reward. An example clarifying the issue is shown in Figure 3.

To allow the agent to use relevant aspects of its history, we will henceforth consider an augmentation of the MDP with auxiliary variables. A theoretical result in the next section states that this choice of augmentation is sufficient in that no additional information will be able to improve performance under the optimal policy.

5.1 AUGMENTATION OF THE RAA PROBLEM

We consider an augmentation of the MDP defined by $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \mathcal{A}, f \rangle$ consisting of augmented states $\overline{\mathcal{S}} = \mathcal{S} \times \mathcal{Y} \times \mathcal{Z}$ and the same actions \mathcal{A} . For any initial state s , let the augmented states be initialized as $y = r(s)$ and $z = q(s)$, and let the transition of $\overline{\mathcal{M}}$ be defined by

$$s_{t+1}^{\bar{\pi}} = f(s_t^{\bar{\pi}}, \bar{\pi}(s_t^{\bar{\pi}}, y_t^{\bar{\pi}}, z_t^{\bar{\pi}})); \quad y_{t+1}^{\bar{\pi}} = \max \{r(s_{t+1}^{\bar{\pi}}), y_t^{\bar{\pi}}\}; \quad z_{t+1}^{\bar{\pi}} = \min \{q(s_{t+1}^{\bar{\pi}}), z_t^{\bar{\pi}}\},$$

such that y_t and z_t track the best reward and worst penalty up to any point. Hence, the policy for $\overline{\mathcal{M}}$ given by $\bar{\pi} : \overline{\mathcal{S}} \rightarrow \mathcal{A}$ may now consider information regarding the history of the trajectory.

By definition, the RAA value for $\overline{\mathcal{M}}$,

$$V_{\text{RAA}}^{\bar{\pi}}(s) = \min \left\{ \max_t r(s_t^{\bar{\pi}}), \min_t q(s_t^{\bar{\pi}}) \right\},$$

is equivalent to that of \mathcal{M} except that it allows for a policy $\bar{\pi}$ which has access to historical information. We seek to find $\bar{\pi}$ that maximizes this value.

5.2 AUGMENTATION OF THE RR PROBLEM

For the Reach-Reach problem, we augment the system similarly, except that z_t is updated using a max operation instead of a min:

$$s_{t+1}^{\bar{\pi}} = f(s_t^{\bar{\pi}}, \bar{\pi}(s_t^{\bar{\pi}}, y_t^{\bar{\pi}}, z_t^{\bar{\pi}})); \quad y_{t+1}^{\bar{\pi}} = \max \{r_1(s_{t+1}^{\bar{\pi}}), y_t^{\bar{\pi}}\}; \quad z_{t+1}^{\bar{\pi}} = \max \{r_2(s_{t+1}^{\bar{\pi}}), z_t^{\bar{\pi}}\}.$$

Again, by definition,

$$V_{\text{RR}}^{\bar{\pi}}(s) = \min \left\{ \max_t r_1(s_t^{\bar{\pi}}), \max_t r_2(s_t^{\bar{\pi}}) \right\}.$$

The RR problem is again to find an augmented policy $\bar{\pi}$ which maximizes this value.

270 6 OPTIMAL POLICIES FOR RAA AND RR BY VALUE DECOMPOSITION

272 We now discuss our first theoretical contributions. We refer the reader to the appendix for the proofs
 273 of the theorems.

275 6.1 DECOMPOSITION OF RAA INTO AVOID AND REACH-AVOID PROBLEMS

277 Our main theoretical result for the RAA problem shows that we can solve this problem by first solving
 278 the avoid problem corresponding to the penalty $q(s)$ to obtain the optimal value function $V_A^*(s)$ and
 279 then solving a reach-avoid problem with the negated penalty function $q(s)$ and a modified reward
 280 function $r_{RAA}(s)$.

281 **Theorem 1.** *For all initial states $s \in \mathcal{S}$,*

$$282 \quad \max_{\bar{\pi}} V_{RAA}^{\bar{\pi}}(s) = \max_{\pi} \max_t \min \left\{ r_{RAA}(s_t^\pi), \min_{\tau \leq t} q(s_\tau^\pi) \right\}, \quad (1)$$

285 where $r_{RAA}(s) := \min \{r(s), V_A^*(s)\}$, with

$$286 \quad V_A^*(s) := \max_{\pi} \min_t q(s_t^\pi).$$

288 This decomposition is significant, as methods customized to solving avoid and reach-avoid problems
 289 were recently explored in Fisac et al. (2019); Hsu et al. (2021); So et al. (2024); So and Fan (2023),
 290 allowing us to effectively solve the optimization problem defining $V_A^*(s)$ as well as the optimization
 291 problem that defines the right-hand-side of 1.

292 **Corollary 1.** *The value function $V_{RAA}^*(s) := \max_{\bar{\pi}} V_{RAA}^{\bar{\pi}}(s)$ satisfies the Bellman equation*

$$294 \quad V_{RAA}^*(s) = \min \left\{ \max \left\{ \max_{a \in \mathcal{A}} V_{RAA}^*(f(s, a)), r_{RAA}(s) \right\}, q(s) \right\}.$$

296 6.2 DECOMPOSITION OF THE RR PROBLEM INTO THREE REACH PROBLEMS

298 Our main result for the RR problem shows that we can solve this problem by first solving two reach
 299 problems corresponding to the rewards $r_1(s)$ and $r_2(s)$ to obtain reach value functions $V_{R1}^*(s)$ and
 300 $V_{R2}^*(s)$, respectively. We then solve a third reach problem with a modified reward $r_{RR}(s)$.

301 **Theorem 2.** *For all initial states $s \in \mathcal{S}$,*

$$303 \quad \max_{\bar{\pi}} V_{RR}^{\bar{\pi}}(s) = \max_{\pi} \max_t r_{RR}(s_t^\pi), \quad (2)$$

304 where $r_{RR}(s) := \max \{\min \{r_1(s), V_{R1}^*(s)\}, \min \{r_2(s), V_{R2}^*(s)\}\}$, with

$$306 \quad V_{R1}^*(s) := \max_{\pi} \max_t r_1(s_t^\pi), \quad V_{R2}^*(s) := \max_{\pi} \max_t r_2(s_t^\pi).$$

308 **Corollary 2.** *The value function $V_{RR}^*(s) := \max_{\bar{\pi}} V_{RR}^{\bar{\pi}}(s)$ satisfies the Bellman equation*

$$309 \quad V_{RR}^*(s) = \max \left\{ \max_{a \in \mathcal{A}} V_{RR}^*(f(s, a)), r_{RR}(s) \right\}.$$

312 6.3 OPTIMALITY OF THE AUGMENTED PROBLEMS

314 We previously motivated the choice to consider an augmented MDP $\overline{\mathcal{M}}$ over the original MDP
 315 in the context of the RAA and RR problems. In this section, we justify our particular choice of
 316 augmentation. Indeed, the following theoretical result shows that further augmenting the states with
 317 additional historical information cannot improve performance under the optimal policy.

318 **Theorem 3.** *Let $s \in \mathcal{S}$. Then*

$$319 \quad \max_{\pi} V_{RAA}^{\pi}(s) \leq \max_{\bar{\pi}} V_{RAA}^{\bar{\pi}}(s) = \max_{a_0, a_1, \dots} \min \left\{ \max_t r(s_t), \min_t q(s_t) \right\},$$

321 and

$$322 \quad \max_{\pi} V_{RR}^{\pi}(s) \leq \max_{\bar{\pi}} V_{RR}^{\bar{\pi}}(s) = \max_{a_0, a_1, \dots} \min \left\{ \max_t r_1(s_t), \max_t r_2(s_t) \right\}$$

323 where $s_{t+1} = f(s_t, a_t)$ and $s_0 = s$.

The terms on the right of the lines above reflect the best possible sequence of actions to solve the RAA or RR problem, and the theorem states that the optimal augmented policy achieves that value, represented by the middle terms. This value will generally be less than or equal to the outcome from using a non-augmented policy, represented by the terms on the left.

7 DO-HJ-PPO: SOLVING RAA AND RR WITH RL

In the previous sections, we demonstrated that the RAA and RR problems can be solved through decomposition of the values into formulations amenable to existing RL methods. However, we make a few assumptions in the derivation that would limit performance and generalization, namely, the determinism of the values as well as access to the decomposed values (by solving them beforehand). In this section, we propose relaxations to the RR and RAA theory and devise a custom variant of Proximal Policy Optimization, **DO-HJ-PPO**, to solve this broader class of problems, and demonstrate its performance.

7.1 STOCHASTIC REACH-AVOID BELLMAN EQUATION

It is well known that the most performative RL methods allow for stochastic learning. In So et al. (2024), the Stochastic Reachability Bellman Equation (SRBE) is described for Reach problems and used to design a specialized PPO algorithm. In this section we proceed by closely following this work, modifying the SRBE into a Stochastic Reach-Avoid Bellman Equation (SRABE). Using Theorems 1 and 2, the SRBE and SRABE offer the necessary tools for designing a PPO variant for solving the RR and RAA problems.

We define $\tilde{V}_{\text{RAA}}^\pi$ to be the solution to the following Bellman equation:

$$\tilde{V}_{\text{RAA}}^\pi(s) = \mathbb{E}_{a \sim \pi} \left[\min \left\{ \max \left\{ \tilde{V}_{\text{RAA}}^\pi(f(s, a)), r_{\text{RAA}}(s) \right\}, q(s) \right\} \right] \quad (\text{SRABE})$$

The corresponding action-value function is

$$\tilde{Q}_{\text{RAA}}^\pi(s, a) = \min \left\{ \max \left\{ \tilde{V}_{\text{RAA}}^\pi(f(s, a)), r_{\text{RAA}}(s) \right\}, q(s) \right\}.$$

We define a modification of the dynamics f involving an absorbing state s_∞ as follows:

$$f'(s, a) = \begin{cases} f(s, a) & q(f(s, a)) < \tilde{V}_{\text{RAA}}^\pi(s) < r_{\text{RAA}}(f(s, a)), \\ s_\infty & \text{otherwise.} \end{cases}$$

We then have the following proposition:

Proposition 1. *For each $s \in \mathcal{S}$ and every $\theta \in \mathbb{R}^{n_p}$, we have*

$$\nabla_\theta \tilde{V}_{\text{RAA}}^{\pi_\theta}(s) \propto \mathbb{E}_{s' \sim d'_\pi(s), a \sim \pi_\theta} \left[\tilde{Q}_{\text{RAA}}^{\pi_\theta}(s', a) \nabla_\theta \ln \pi_\theta(a|s') \right],$$

where $d'_\pi(s)$ is the stationary distribution of the Markov Chain with transition function

$$P(s'|s) = \sum_{a \in \mathcal{A}} \pi(a|s) [f'(s, \pi(a|s)) = s'],$$

with the bracketed term equal to 1 if the proposition inside is true and 0 otherwise.

Following Hsu et al. (2021), we then define the discounted value and action-value functions with $\gamma \in [0, 1]$.

$$\tilde{V}_{\text{RAA}}^{\gamma, \pi}(s) = (1 - \gamma) \min \{r_{\text{RAA}}(s), q(s)\} + \gamma \mathbb{E}_{a \sim \pi} \left[\min \left\{ \max \left\{ \tilde{V}_{\text{RAA}}^{\gamma, \pi}(f(s, a)), r_{\text{RAA}}(s) \right\}, q(s) \right\} \right].$$

$$\tilde{Q}_{\text{RAA}}^{\gamma, \pi}(s, a) = (1 - \gamma) \min \{r_{\text{RAA}}(s), q(s)\} + \gamma \min \left\{ \max \left\{ \tilde{V}_{\text{RAA}}^{\gamma, \pi}(f(s, a)), r_{\text{RAA}}(s) \right\}, q(s) \right\}.$$

The PPO advantage function is then given by $\tilde{A}_{\text{RAA}}^\pi = \tilde{Q}_{\text{RAA}}^\pi - \tilde{V}_{\text{RAA}}^\pi$ Schulman et al. (2017).

378
379

7.2 ALGORITHM

380
381

We introduce **DO-HJ-PPO** for solving the RAA and RR problems, which integrates the SRABE and SRBE via three minimal modifications to PPO Schulman et al. (2017) (see appendix for more).

382
383
384
385
386

Additional actor and critics are introduced to represent the decomposed objectives. Per Theorems 1 and 2, one may know that the RAA and RR values are given by a composition of the simpler R, A and RA values. Therefore, we learn these decompositions with their own networks and integrate them into the composed actor and critic training, namely via the GAE and target with the special RAA and RR reward functions in Theorems 1 and 2.

387
388
389
390
391
392
393

The composed actor and critic are learned concurrently to the decomposed actor and critics by bootstrapping the current values Rather than learning the decomposed and composed representations sequentially, DO-HJ-PPO bootstraps to learn them simultaneously. Namely, at each iteration, we rollout trajectories for composed and decomposed updates with each actor. In the update of the composed representation specifically, the decomposed values are inferred from the current decomposed critic(s) along the composed trajectories. This design choice allows us to couple the on-policy learning of PPO in the following way.

394
395
396
397
398
399
400

Trajectories for training the decomposed actor and critic(s) are initialized with states sampled from the composed trajectories, which we refer to as *coupled resets*. While it is possible to estimate the decomposed objectives independently—i.e., prior to solving the composed task—this approach might lead to inaccurate or irrelevant value estimates in on-policy settings. For example, in the RAA problem, the avoid decomposition will solely prioritize avoiding penalties and, hence, might converge to an optimal strategy within a reward-irrelevant region, misaligned with the overall task.

401
402

8 EXPERIMENTS

403
404

8.1 DQN DEMONSTRATION

405
406
407
408
409
410
411
412
413
414
415
416
417
418
419

We begin by demonstrating the utility of our theoretical results (Theorems 1 and 2) through a simple 2D grid-world experiment using DQN (Figure 2). In this environment, the agent can move left, right, or remain stationary, while drifting upward at a constant rate. Throughout, reward regions are shown in blue and penalty regions in red. On the left, we compare the optimal value functions learned under the classic Reach-Avoid (RA) formulation with those from the Reach-Always-Avoid (RAA) setting. In the RA scenario, trajectories successfully avoid the obstacle but may terminate in regions from which future collisions are inevitable, as there is no incentive to consider what happens after reaching the minimum reward threshold. In contrast, under the RAA formulation, where the objective involves maximizing cumulative reward while accounting for future penalties (as per Theorem 1), the agent learns to reach the target while remaining in safe regions thereafter. On the right, we consider a similar environment without obstacles but with two distinct targets. Here, the Reach-Reach (RR) formulation induces trajectories that visit both targets, unlike simple reach tasks in which the agent halts after reaching a single goal. These qualitative results highlight the behavioral distinctions induced by the RAA and RR objectives compared to their simpler counterparts. Additional algorithmic and experimental details are provided in the Appendix.

420
4218.2 CONTINUOUS CONTROL TASKS WITH **DO-HJ-PPO**422
423
424
425
426

To evaluate the method under more complex and less structured conditions, we extend our analysis to continuous control settings. Specifically, we consider RAA and RR tasks in the Hopper, F16, SafetyGym, and HalfCheetah environments, depicted in Figure 1. In the RAA tasks, the penalty function generally characterizes regions of states where the agent (or its body parts) is intended to avoid, while the reward characterizes regions of states where the agent is intended to reach.

427
428
429
430
431

As baselines, we compare **DO-HJ-PPO** against a variety of classes of RL algorithms. We include several augmented Lagrangian methods which transform constraints (either for reaching both or always avoiding) into mixed objectives, such as Constrained PPO (CPPO) Achiam et al. (2017b), PPO-LAG Ray et al. (2019), P2BPO Dey et al. (2024), and LOGBAR Zhang et al. (2024). Additionally, we include three HJ-RL baselines designed for the previous R and RA problems, RESPO Ganai et al. (2023), RCPPO So et al. (2024) and RA Hsu et al. (2021). Lastly, we also include a few methods

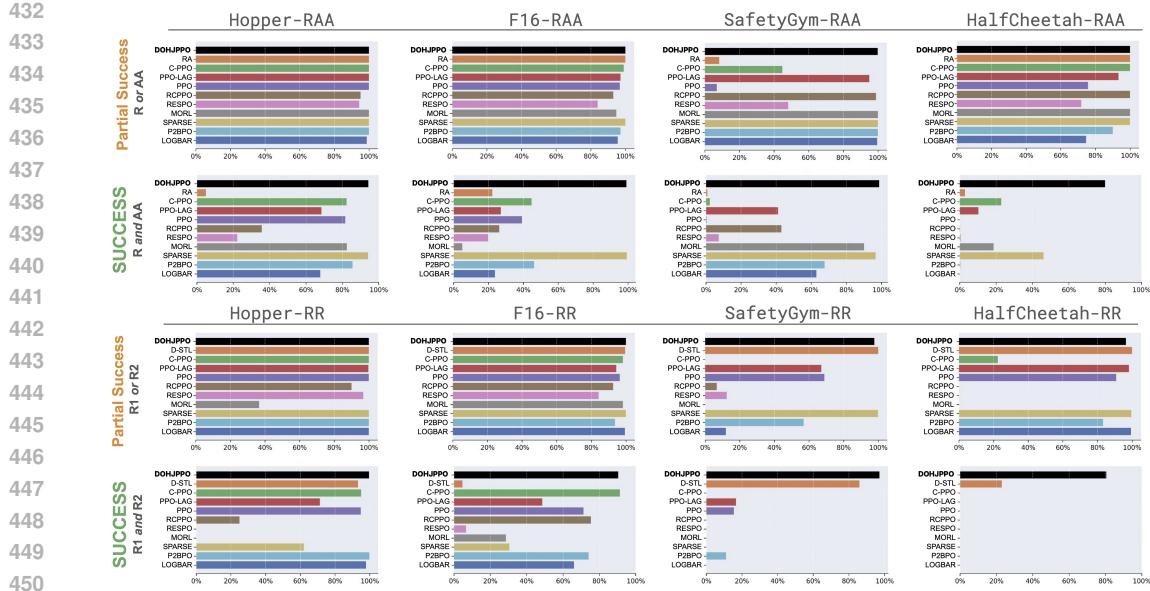


Figure 4: Success (\rightarrow) and Partial Success (\rightarrow) in RAA and RR Tasks for DO-HJ-PPO and Baselines
 We evaluate DO-HJ-PPO in black against baselines over 1,000 trajectories in the Hopper, F16, SafetyGym and HalfCheetah environments. In the first and third row, the **partial success** percentage of each algorithm is given, defined by the number of trajectories to achieve one objective (reaching or always-avoiding in the RAA, reaching either in the RR). In the second and fourth rows, **success** percentage is given, defined by the number of trajectories to achieve both objectives. Most baselines achieve partial success, however, few achieve total success as the environment becomes more difficult, underscoring the difficulty of balancing objectives in RL.

based on approaches in STL/LTL-RL and MORL, including a decomposed STL (D-STL) PPO, a sparse-reward STL PPO (SPARSE) and a MORL-based PPO. All algorithms are trained on random initial conditions and then evaluated on new random initial conditions within distribution. To quantify performance of the dual-objective tasks, we measure (1) the percent of trajectories which achieve at least both tasks successfully, (2) the percent of trajectories which achieve at least one of the tasks (dubbed partial success), and (3) the mean steps in each trajectory until success.

Empirically, we find that our method performs at the top-level, achieving first or second place among all tasks and environments (Figure 4). In fact, as the multi-target (RR) or safe-achievement (RAA) tasks become more complex (e.g. the HalfCheetah), our algorithm increasingly dominates the 10 state-of-the-art baselines. Note, that almost all algorithms can achieve partial success at a high rate in each dual-objective task, highlighting the difficulty of mixed or competing objectives, particularly with discounted-sum rewards. Moreover, DO-HJ-PPO is the sole performant algorithm in both RAA and RR tasks, displaying the fastest achievement times across tasks (see appendix).

These results underscore the challenging nature of composing multiple satisfaction objectives using traditional baselines with discounted-sum rewards. In contrast, DO-HJ-PPO provides a direct and robust solution to handling these complex tasks, with *little to no* tuning. Our algorithm enjoys these benefits because of the structure of the novel Bellman updates, which propagate the extreme (maximum and minimum) values as opposed to the short-term average (discounted-sum) values.

9 CONCLUSION

In this work, we introduced two novel Bellman formulations for new problems (RAA and RR) which generalize those considered in several recent publications. We derive decomposition results to break them into simpler Bellman equations, which can then be composed to obtain the corresponding value functions and optimal policies. We use these results to design DO-HJ-PPO for practical solution of RAA and RR, which proves the most performant and balanced algorithm in safe-arrival and multi-target achievement. More broadly, this work provides a road-map to extend complex Bellman formulations, via decomposing higher-level problems into lower-level ones, establishing a foundation for nuanced tasks in real-world environments and safe RL.

486 REFERENCES
487

- 488 Axel Abels, Diederik Roijers, Tom Lenaerts, Ann Nowé, and Denis Steckelmacher. Dynamic weights
489 in multi-objective deep reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov,
490 editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
491 *Proceedings of Machine Learning Research*, pages 11–20. PMLR, 2019.
- 492 Joshua Achiam, David Held, Aviv Tamar, and P Abbeel. Constrained policy optimization. *ICML*,
493 abs/1705.10528:22–31, 30 May 2017a.
- 494 Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In
495 *International conference on machine learning*, pages 22–31. PMLR, 2017b.
- 496 Eitan Altman. *Constrained Markov decision processes: Stochastic modeling*. Routledge, Boca Raton,
497 13 December 2021.
- 498 Aaron D Ames, Xiangru Xu, Jessy W Grizzle, and Paulo Tabuada. Control barrier function based
499 quadratic programs for safety critical systems. *IEEE Transactions on Automatic Control*, 62(8):
500 3861–3876, 2016.
- 501 F Bacchus, Craig Boutilier, and Adam J Grove. Rewarding behaviors. In *Proceedings of the National*
502 *Conference on Artificial Intelligence.*, pages 1160–1167. cs.toronto.edu, 4 August 1996.
- 503 Fahiem Bacchus, Craig Boutilier, and Adam Grove. Structured solution methods for non-Markovian
504 decision processes. In *AAAI/IAAI*, pages 112–117, 1997.
- 505 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
506 Wojciech Zaremba. Openai gym, 2016. URL <https://arxiv.org/abs/1606.01540>.
- 507 Tomáš Brázdil, Krishnendu Chatterjee, Martin Chmelík, Vojtěch Forejt, Jan Křetínský, Marta
508 Kwiatkowska, David Parker, and Mateusz Ujma. Verification of Markov decision processes
509 using learning algorithms. *arXiv [cs.LO]*, 10 February 2014.
- 510 Xin-Qiang Cai, Pushi Zhang, Li Zhao, Jiang Bian, Masashi Sugiyama, and Ashley Llorens. Dis-
511 tributional Pareto-optimal multi-objective reinforcement learning. *Neural Inf Process Syst*, 36:
512 15593–15613, 2023.
- 513 Alberto Camacho, Rodrigo Toro Icarte, Toryn Q Klassen, Richard Valenzano, and Sheila A McIlraith.
514 LTL and beyond: Formal languages for reward function specification in reinforcement learning. In
515 *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, pages
516 6065–6073, California, 1 August 2019. International Joint Conferences on Artificial Intelligence
517 Organization.
- 518 Alberto Camacho, Oscar Chen, Scott Sanner, and Sheila McIlraith. Non-Markovian rewards expressed
519 in LTL: Guiding search via reward shaping. *Proceedings of the International Symposium on*
520 *Combinatorial Search*, 8(1):159–160, 1 September 2021.
- 521 Andres Campero, Roberta Raileanu, Heinrich Küttler, Joshua B Tenenbaum, Tim Rocktäschel, and
522 Edward Grefenstette. Learning with AMIGo: Adversarially motivated intrinsic goals. *arXiv*
523 [*cs.LG*], 22 June 2020.
- 524 Elliot Chane-Sane, C Schmid, and I Laptev. Goal-conditioned reinforcement learning with imagined
525 subgoals. *ICML*, abs/2107.00541:1430–1440, 1 July 2021.
- 526 Yi Chen, Jing Dong, and Zhaoran Wang. A primal-dual approach to constrained Markov decision
527 processes. *arXiv [math.OC]*, 26 January 2021.
- 528 Max H Cohen, Zachary Serlin, Kevin Leahy, and Calin Belta. Temporal logic guided safe model-
529 based reinforcement learning: A hybrid systems approach. *Nonlinear Anal. Hybrid Syst.*, 47
530 (101295):101295, February 2023.
- 531 Sumanta Dey, Pallab Dasgupta, and Soumyajit Dey. P2bpo: Permeable penalty barrier-based
532 policy optimization for safe rl. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
533 volume 38, pages 21029–21036, 2024.

- 540 Benjamin Eysenbach, Tianjun Zhang, R Salakhutdinov, and S Levine. Contrastive learning as
 541 goal-conditioned reinforcement learning. *Neural Inf Process Syst*, abs/2206.07568:35603–35620,
 542 15 June 2022.
- 543 Jaime F Fisac, Mo Chen, Claire J Tomlin, and S Shankar Sastry. Reach-avoid problems with time-
 544 varying dynamics, targets and constraints. In *Hybrid Systems: Computation and Control*. ACM,
 545 2015.
- 546 Jaime F Fisac, Neil F Lugovoy, Vicenç Rubies-Royo, Shromona Ghosh, and Claire J Tomlin. Bridging
 547 hamilton-jacobi safety analysis and reinforcement learning. In *2019 International Conference on*
 548 *Robotics and Automation (ICRA)*, pages 8550–8556. IEEE, 2019.
- 549 Milan Ganai, Chiaki Hirayama, Ya-Chien Chang, and Sicun Gao. Learning stabilization control from
 550 observations by learning lyapunov-like proxy models. *2023 IEEE International Conference on*
 551 *Robotics and Automation (ICRA)*, 2023.
- 552 Nathaniel Hamilton, Preston K Robinette, and Taylor T Johnson. Training agents to satisfy timed and
 553 untimed signal temporal logic specifications with reinforcement learning. In *Software Engineering*
 554 and *Formal Methods*, Lecture notes in computer science, pages 190–206. Springer International
 555 Publishing, Cham, 2022.
- 556 Kai-Chieh Hsu, Vicenç Rubies-Royo, Claire J. Tomlin, and Jaime F. Fisac. Safety and liveness
 557 guarantees through reach-avoid reinforcement learning. In *Proceedings of Robotics: Science and*
 558 *Systems*, Held Virtually, July 2021. doi: 10.15607/RSS.2021.XVII.077.
- 559 Rodrigo Toro Icarte, Toryn Q Klassen, R Valenzano, and Sheila A McIlraith. Using reward machines
 560 for high-level task specification and decomposition in reinforcement learning. *ICML*, 80:2112–
 561 2121, 3 July 2018.
- 562 Tianjiao Li, Ziwei Guan, Shaofeng Zou, Tengyu Xu, Yingbin Liang, and Guanghui Lan. Faster
 563 algorithm and sharper analysis for constrained Markov decision process. *Oper. Res. Lett.*, 54
 564 (107107):107107, May 2024.
- 565 Xiao Li, Cristian-Ioan Vasile, and Calin Belta. Reinforcement learning with temporal logic rewards.
 566 In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
 567 3834–3839. IEEE, September 2017.
- 568 Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning: Problems
 569 and solutions. *arXiv [cs.AI]*, 20 January 2022.
- 570 Jason Yecheng Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. Offline goal-conditioned rein-
 571 forcement learning via f-advantage regression. In S Koyejo, S Mohamed, A Agarwal, D Belgrave,
 572 K Cho, and A Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages
 573 310–323. Curran Associates, Inc., 2022.
- 574 Sobhan Miryoosefi and Chi Jin. A simple reward-free approach to constrained reinforcement learning.
 575 *ICML*, abs/2107.05216:15666–15698, 12 July 2021.
- 576 Ian M Mitchell, Alexandre M Bayen, and Claire J Tomlin. A time-dependent hamilton-jacobi
 577 formulation of reachable sets for continuous dynamic games. *IEEE Transactions on automatic*
 578 *control*, 50(7):947–957, 2005.
- 579 Hossam Mossalam, Yannis M Assael, Diederik M Roijers, and Shimon Whiteson. Multi-objective
 580 deep reinforcement learning. *arXiv [cs.AI]*, 9 October 2016.
- 581 Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Powell,
 582 Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, Vikash Kumar, and Wojciech
 583 Zaremba. Multi-goal reinforcement learning: Challenging robotics environments and request for
 584 research. *arXiv [cs.LG]*, 26 February 2018.
- 585 Wenjie Qiu, Wensen Mao, and He Zhu. Instructing goal-conditioned reinforcement learning agents
 586 with temporal logic objectives. *Neural Inf Process Syst*, 36:39147–39175, 2023.

- 594 Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement
 595 learning. *arXiv preprint arXiv:1910.01708*, 7(1):2, 2019.
- 596
- 597 Zhizhou Ren, Kefan Dong, Yuanshuo Zhou, Qiang Liu, and Jian Peng. Exploration via hindsight
 598 goal generation. *Neural Inf Process Syst*, 32:13464–13474, 1 June 2019.
- 599
- 600 Dorsa Sadigh, Eric S Kim, Samuel Coogan, S Shankar Sastry, and Sanjit A Seshia. A learning based
 601 approach to control synthesis of Markov decision processes for linear temporal logic specifications.
 In *53rd IEEE Conference on Decision and Control*, pages 1091–1096. IEEE, December 2014.
- 602
- 603 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
 604 optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- 605
- 606 Oswin So and Chuchu Fan. Solving stabilize-avoid optimal control via epigraph form and deep
 607 reinforcement learning. *arXiv [cs.RO]*, 23 May 2023.
- 608
- 609 Oswin So, Cheng Ge, and Chuchu Fan. Solving minimum-cost reach avoid using reinforcement
 610 learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,
 2024. URL <https://openreview.net/forum?id=jzngdJQ21Y>.
- 611
- 612 Adam Stooke, Joshua Achiam, and P Abbeel. Responsive safety in reinforcement learning by PID
 613 lagrangian methods. *ICML*, 119:9133–9143, 8 July 2020.
- 614
- 615 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford
 Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- 616
- 617 S Thiebaux, C Gretton, J Slaney, D Price, and F Kabanza. Decision-theoretic planning with non-
 618 Markovian rewards. *J. Artif. Intell. Res.*, 25:17–74, 29 January 2006.
- 619
- 620 Alexander R Trott, Stephan Zheng, Caiming Xiong, and R Socher. Keeping your distance: Solving
 621 sparse reward tasks using self-balancing shaped rewards. *Neural Inf Process Syst*, abs/1911.01417,
 4 November 2019.
- 622
- 623 Moffaert K Van and A Nowé. Multi-objective reinforcement learning using sets of Pareto dominating
 policies. *The Journal of Machine Learning Research*, 15(1):3483–3512, 2014.
- 624
- 625 Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-
 626 learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- 627
- 628 Akifumi Wachi and Yanan Sui. Safe reinforcement learning in constrained Markov decision processes.
ICML, 119:9797–9806, 12 July 2020.
- 629
- 630 Marco A Wiering, Maikel Withagen, and Madalina M Drugan. Model-based multi-objective reinforce-
 631 ment learning. In *2014 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement
 Learning (ADPRL)*, pages 1–6. IEEE, December 2014.
- 632
- 633 R Yang, X Sun, and K Narasimhan. A generalized algorithm for multi-objective reinforcement
 634 learning and policy adaptation. In *Advances in Neural Information Processing Systems*. proceedings.neurips.cc, 2019.
- 635
- 636 Tsung-Yen Yang, Justinian Rosca, Karthik Narasimhan, and Peter J Ramadge. Projection-based
 637 constrained policy optimization. *arXiv [cs.LG]*, 7 October 2020.
- 638
- 639 Dongjie Yu, Haitong Ma, Shengbo Li, and Jianyu Chen. Reachability constrained reinforcement
 640 learning. In *International Conference on Machine Learning*, pages 25636–25655. PMLR, 2022a.
- 641
- 642 Dongjie Yu, Wenjun Zou, Yujie Yang, Haitong Ma, Shengbo Eben Li, Jingliang Duan, and Jianyu
 643 Chen. Safe model-based reinforcement learning with an uncertainty-aware reachability certificate.
arXiv preprint arXiv:2210.07553, 2022b.
- 644
- 645 Baohe Zhang, Yuan Zhang, Lilli Frison, Thomas Brox, and Joschka Bödecker. Constrained reinforce-
 646 ment learning with smoothed log barrier function. *arXiv preprint arXiv:2403.14508*, 2024.
- 647
- Kai Zhu, Fengbo Lan, Wenbo Zhao, and Tao Zhang. Safe multi-agent reinforcement learning via
 approximate hamilton-jacobi reachability. *J. Intell. Robot. Syst.*, 111(1), 30 December 2024.

648
649 **10 ETHICS STATEMENT**

650 This project was conducted and completed on entirely responsible and ethical grounds, and meets the
651 highest standard of the ICLR Code of Ethics. We believe the rigor, investigation and communication
652 not only upholds scientific ideals whilst avoiding societal harm, but advances machine learning for
653 the betterment of all society, namely by improving learning to be more performant with much less
654 hyper-tuning, and thus better for the planet and human race. Moreover, the work fundamentally
655 improves the safety and reliability of reinforcement learning, and thus greatly improves a society in
656 which machine learning is heavily integrated. Above all, the work is honest, noting limitations and
657 caveats, while depicting the strengths we believe make this work invaluable for the field.

658

659 **11 REPRODUCIBILITY STATEMENT**

660

661 All theorems, algorithms and parameters for this work are totally explained in this paper (partially in
662 the appendix) in what the authors believe to be a clear and understandable form. All theorems have
663 been proven in detail, including all necessary lemmas, propositions and references in a manner the
664 authors believe is intelligible. The algorithm proposed in the work is explained clearly in the main
665 text and written line-by-line in the appendix, along with all hyper-parameters for each environment.
666 The code for this work has been inherited from another group with security clearance and we are
667 awaiting their response to publicize it, but will do so as soon as possible as we are committed to fair
668 and open resources without bias or discrimination.

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

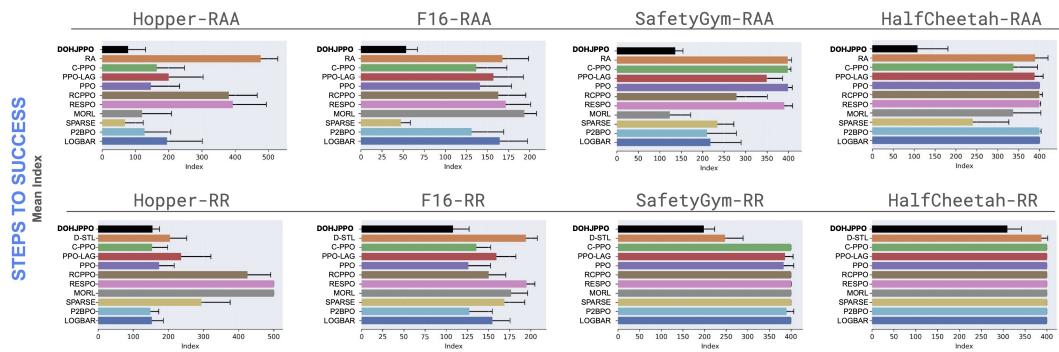
702 Appendix

703 Contents

707 α Achievement Speed Results from DO-HJ-PPO Experiments	14
708 A Proof of RAA Main Theorem	15
709 B Proof of RR Main Theorem	23
710 C Proof of Optimality Theorem	28
711 D The SRABE and its Policy Gradient	28
712 E The DO-HJ-PPO Algorithm	29
713 F DDQN Demonstration	30
714 G Baselines	31
715 H Details of RAA & RR Experiments: Hopper	32
716 I Details of RAA & RR Experiments: F16	33
717 J Broader Impacts	34
718 K Acknowledgments	34

729 ACHIEVEMENT SPEED RESULTS FROM DO-HJ-PPO EXPERIMENTS

730 Here we present additional results for RAA and RR problems solved with **DO-HJ-PPO**. In both
 731 settings, DO-HJ-PPO out-performs or matches the best of baselines with less tuning and faster arrival.
 732 Notably as the difficulty of the problem increases the gap increases significantly with DO-HJ-PPO
 733 remaining the sole algorithm that can achieve the task in reasonable time and in both RAA and RR
 734 categories.



748 **Figure 5: Steps to Success (←) in RAA and RR Tasks for DO-HJ-PPO and Baselines** For the same 1000
 749 trajectories in Figure 4, we quantify here the number of steps until achievement of both tasks: reaching without
 750 crash afterward in the RAA, reaching both goal in the RR. DO-HJ-PPO is not only competitive but consistently
 751 achieves the dual-objective problems in the fewest number of steps.

752 PROOF NOTATION

753 Throughout the theoretical sections of this supplement, we use the following notation.

754 We let $\mathbb{N} = \{0, 1, \dots\}$ be the set of whole numbers.

We let \mathbb{A} be the set of maps from \mathbb{N} to \mathcal{A} . In other words, \mathbb{A} is the set of sequences of actions the agent can choose. Given $\mathbf{a}_1, \mathbf{a}_2 \in \mathbb{A}$, and $\tau \in \mathbb{N}$, we let $[\mathbf{a}_1, \mathbf{a}_2]_\tau$ be the element of \mathbb{A} for which

$$[\mathbf{a}_1, \mathbf{a}_2]_\tau(t) = \begin{cases} \mathbf{a}_1(t) & t < \tau, \\ \mathbf{a}_2(t - \tau) & t \geq \tau. \end{cases}$$

Similarly, given $a \in \mathcal{A}$ and $\mathbf{a} \in \mathbb{A}$, we let $[a, \mathbf{a}]$ be the element of \mathbb{A} for which

$$[a, \mathbf{a}](t) = \begin{cases} a & t = 0, \\ \mathbf{a}(t - 1) & t \geq 1. \end{cases}$$

Additionally, given $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$, we let $\mathbf{a}|_\tau$ be the element of \mathbb{A} for which

$$\mathbf{a}|_\tau(t) = \mathbf{a}(t + \tau) \quad \forall t \in \mathbb{N}.$$

The $[\cdot, \cdot]_\tau$ operation corresponds to concatenating two action sequences (using only the 0th to $(\tau - 1)$ st elements of the first sequence), the $[\cdot, \cdot]$ operation corresponds to prepending an action to an action sequence, and the $\cdot|_\tau$ operation corresponds to removing the 0th to $(\tau - 1)$ st elements of an action sequence.

We let Π be the set of policies $\pi : \mathcal{S} \rightarrow \mathcal{A}$. Given $s \in \mathcal{S}$ and $\pi \in \Pi$, we let $\xi_s^\pi : \mathbb{N} \rightarrow \mathcal{S}$ be the solution of the evolution equation

$$\xi_s^\pi(t + 1) = f(\xi_s^\pi(t), \pi(\xi_s^\pi(t)))$$

for which $\xi_s^\pi(0) = s$. In other words, $\xi_s^\pi(\cdot)$ is the state trajectory over time when the agent begins at state s and follows policy π .

We will also “overload” this trajectory notation for signals rather than policies: given $\mathbf{a} \in \mathbb{A}$, we let $\xi_s^\mathbf{a} : \mathbb{N} \rightarrow \mathcal{S}$ be the solution of the evolution equation

$$\xi_s^\mathbf{a}(t + 1) = f(\xi_s^\mathbf{a}(t), \mathbf{a}(t))$$

for which $\xi_s^\mathbf{a}(0) = s$. In other words, $\xi_s^\mathbf{a}(\cdot)$ is the state trajectory over time when the agent begins at state s and follows action sequence \mathbf{a} .

A PROOF OF RAA MAIN THEOREM

We first define the value functions, $V_A^*, \tilde{V}_{RA}^*, V_{RAA}^* : \mathcal{S} \rightarrow \mathbb{R}$ by

$$\begin{aligned} V_A^*(s) &= \max_{\pi \in \Pi} \min_{\tau \in \mathbb{N}} q(\xi_s^\pi(\tau)), \\ \tilde{V}_{RA}^*(s) &= \max_{\pi \in \Pi} \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\}, \\ V_{RAA}^*(s) &= \max_{\pi \in \Pi} \min \left\{ \max_{\tau \in \mathbb{N}} r(\xi_s^\pi(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^\pi(\kappa)) \right\}, \end{aligned}$$

where r_{RAA} is as in Theorem 1.

We next define the value functions, $v_A^*, \tilde{v}_{RA}^*, v_{RAA}^* : \mathcal{S} \rightarrow \mathbb{R}$, which maximize over action sequences rather than policies:

$$\begin{aligned} v_A^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \min_{\tau \in \mathbb{N}} q(\xi_s^\mathbf{a}(\tau)), \\ \tilde{v}_{RA}^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^\mathbf{a}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\mathbf{a}(\kappa)) \right\}, \\ v_{RAA}^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \min \left\{ \max_{\tau \in \mathbb{N}} r(\xi_s^\mathbf{a}(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^\mathbf{a}(\kappa)) \right\}, \end{aligned}$$

Observe that for each $s \in \mathcal{S}$,

$$v_A^*(s) \geq V_A^*(s), \quad \tilde{v}_{RA}^*(s) \geq \tilde{V}_{RA}^*(s), \quad v_{RAA}^*(s) \geq V_{RAA}^*(s).$$

We now prove a series of lemmas that will be useful in the proof of the main theorem.

810 **Lemma 1.** *There is a $\pi \in \Pi$ such that*

$$812 \quad v_A^*(s) = \min_{\tau \in \mathbb{N}} q(\xi_s^\pi(\tau))$$

813 *for all $s \in \mathcal{S}$.*

815 *Proof.* Choose $\pi \in \Pi$ such that

$$817 \quad \pi(s) \in \arg \max_{a \in \mathcal{A}} v_A^*(f(s, a)) \quad \forall s \in \mathcal{S}.$$

819 Fix $s \in \mathcal{S}$. Note that for each $\tau \in \mathbb{N}$,

$$\begin{aligned} 820 \quad v_A^*(\xi_s^\pi(\tau + 1)) &= v_A^*(f(\xi_s^\pi(\tau), \pi(\xi_s^\pi(\tau)))) \\ 821 &= \max_{a \in \mathcal{A}} v_A^*(f(\xi_s^\pi(\tau), a)) \\ 822 &= \max_{a \in \mathcal{A}} \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q\left(\xi_{f(\xi_s^\pi(\tau), a)}^{\mathbf{a}}(\kappa)\right) \\ 823 &= \max_{a \in \mathcal{A}} \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q\left(\xi_{\xi_s^\pi(\tau)}^{[a, \mathbf{a}]}(\kappa + 1)\right) \\ 824 &= \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q\left(\xi_{\xi_s^\pi(\tau)}^{\mathbf{a}}(\kappa + 1)\right) \\ 825 &\geq \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q\left(\xi_{\xi_s^\pi(\tau)}^{\mathbf{a}}(\kappa)\right) \\ 826 &\geq v_A^*(\xi_s^\pi(\tau)). \end{aligned}$$

827 It follows by induction that $v_A^*(\xi_s^\pi(\tau)) \geq v_A^*(\xi_s^\pi(0))$ for all $\tau \in \mathbb{N}$, so that

$$828 \quad v_A^*(s) \geq \min_{\tau \in \mathbb{N}} q(\xi_s^\pi(\tau)) \geq \min_{\tau \in \mathbb{N}} v_A^*(\xi_s^\pi(\tau)) = v_A^*(\xi_s^\pi(0)) = v_A^*(s).$$

830 \square

831 **Corollary 3.** *For all $s \in \mathcal{S}$, we have $V_A^*(s) = v_A^*(s)$.*

832 **Lemma 2.** *There is a $\pi \in \Pi$ such that*

$$833 \quad \tilde{v}_{RA}^*(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\}$$

834 *for all $s \in \mathcal{S}$.*

835 *Proof.* First, let us note that in this proof we will use the standard conventions that

$$836 \quad \max \emptyset = -\infty \quad \text{and} \quad \min \emptyset = +\infty.$$

837 We next introduce some notation. First, for convenience, we set $v^* = \tilde{v}_{RA}^*$ and $V^* = \tilde{V}_{RA}^*$. Given
838 $s \in \mathcal{S}$ and $\mathbf{a} \in \mathbb{A}$, we write

$$839 \quad v^{\mathbf{a}}(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}.$$

840 Similarly, given $s \in \mathcal{S}$ and $\pi \in \Pi$, we write

$$841 \quad V^\pi(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\}.$$

842 Then

$$843 \quad V^*(s) = \max_{\pi \in \Pi} \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\} = \max_{\pi \in \Pi} V^\pi(s),$$

844 and

$$845 \quad v^*(s) = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r_{RAA}(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} = \max_{\mathbf{a} \in \mathbb{A}} v^{\mathbf{a}}(s).$$

846 It is immediate that $v^*(s) \geq V^*(s)$ for each $s \in \mathcal{S}$, so it suffices to show the reverse inequality.
847 Toward this end, it suffices to show that there is a $\pi \in \Pi$ for which $V^\pi(s) = v^*(s)$ for each $s \in \mathcal{S}$.
848 Indeed, in this case, $V^*(s) \geq V^\pi(s) = v^*(s)$.

We now construct the desired policy π . Let $\alpha_0 = +\infty$, $S_0 = \emptyset$, and $v_0^* : \mathcal{S} \rightarrow \mathbb{R} \cup \{-\infty\}$, $s \mapsto -\infty$. We recursively define $\alpha_t \in \mathbb{R}$, $S_t \subseteq \mathcal{S}$, and $v_t^* : \mathcal{S} \rightarrow \mathbb{R} \cup \{-\infty\}$ for $t = 1, 2, \dots$ by

$$\alpha_{t+1} = \max_{s \in \mathcal{S} \setminus S_t} \min \left\{ \max \left\{ r_{\text{RAA}}(s), \max_{a \in \mathcal{A}} v_t^*(f(s, a)) \right\}, q(s) \right\}, \quad (3)$$

$$S_{t+1} = S_t \cup \left\{ s \in \mathcal{S} \setminus S_t \mid \min \left\{ \max \left\{ r_{\text{RAA}}(s), \max_{a \in \mathcal{A}} v_t^*(f(s, a)) \right\}, q(s) \right\} = \alpha_{t+1} \right\}, \quad (4)$$

$$v_{t+1}^*(s) = \begin{cases} v_t^*(s) & s \in S_t, \\ \alpha_{t+1} & s \in S_{t+1} \setminus S_t, \\ -\infty & s \in \mathcal{S} \setminus S_{t+1}. \end{cases} \quad (5)$$

From (4) it follows that

$$S_0 \subseteq S_1 \subseteq S_2 \subseteq \dots, \quad (6)$$

which together with (3) shows that

$$\alpha_0 \geq \alpha_1 \geq \alpha_2 \geq \dots. \quad (7)$$

Also, whenever $\mathcal{S} \setminus S_t$ is non-empty, the set being appended to S_t in (4) is non-empty so

$$\bigcup_{t=0}^{\infty} S_t = \mathcal{S}. \quad (8)$$

For each $s \in \mathcal{S}$, let $\sigma(s)$ be the smallest $t \in \mathbb{N}$ for which $s \in S_t$. We choose the policy $\pi \in \Pi$ of interest by insisting

$$\pi(s) \in \arg \max_{a \in \mathcal{A}} v_{\sigma(s)-1}^*(f(s, a)) \quad \forall s \in \mathcal{S}. \quad (9)$$

In the remainder of the proof, we show that $V^\pi(s) = v^*(s)$ for each $s \in \mathcal{S}$ by induction. Let $n \in \mathbb{N}$ and suppose the following induction assumptions hold:

$$V^\pi(s) = v^*(s) = v_n^*(s) \geq \alpha_n \quad \forall s \in S_n, \quad (10)$$

$$v^*(s') \leq \alpha_n \quad \forall s' \in \mathcal{S} \setminus S_n. \quad (11)$$

Note that the above hold trivially when $n = 0$ since $S_0 = \emptyset$ and $\alpha_0 = +\infty$. Fix some particular $y \in S_{n+1}$ and some $z \in \mathcal{S} \setminus S_{n+1}$. We must show that

$$V^\pi(y) = v^*(y) = v_{n+1}^*(y) \geq \alpha_{n+1}, \quad (12)$$

$$v^*(z) \leq \alpha_{n+1}. \quad (13)$$

In this case, induction then shows that $V^\pi(s) = v^*(s)$ for all $s \in \bigcup_{t=0}^{\infty} S_t$. Since this union is equal to \mathcal{S} by (8), the desired result then follows.

To show (12)-(13), we first demonstrate the following three claims.

- Let $x \in \mathcal{S}$ and $w \in \mathcal{A}$ be such that $f(x, w) \in S_n$ and $q(x) \geq \alpha_{n+1}$. We claim $x \in S_{n+1}$.

We can assume $x \notin S_n$, for otherwise the claim follows immediately from (6). Since $f(x, w) \in S_n$, we have $v_n^*(f(x, w)) \geq \alpha_n$ by (10). Thus

$$\begin{aligned} \alpha_{n+1} &\geq \min \left\{ \max \left\{ r_{\text{RAA}}(x), \max_{a \in \mathcal{A}} v_n^*(f(x, a)) \right\}, q(x) \right\} \\ &\geq \min \{ \max \{ r_{\text{RAA}}(x), \alpha_n \}, \alpha_{n+1} \} \\ &= \alpha_{n+1}, \end{aligned}$$

where the first inequality follows from (3), and the equality follows from (7). Thus

$$\alpha_{n+1} = \min \left\{ \max \left\{ r_{\text{RAA}}(x), \max_{a \in \mathcal{A}} v_n^*(f(x, a)) \right\}, q(x) \right\},$$

so the claim follows from (4).

- 918 2. Let $x \in S_{n+1} \setminus S_n$ and $w \in \mathcal{A}$ be such that $f(x, w) \in S_n$. We claim that
 919
 920

$$V^\pi(x) = v^*(x) = \alpha_{n+1}. \quad (14)$$

921 To show this claim, we will make use of the dynamic programming principle
 922
 923

$$v^{\mathbf{a}}(s) = \min \left\{ \max \left\{ r_{\text{RAA}}(s), v^{\mathbf{a}|_1}(f(s, \mathbf{a}(0))) \right\}, q(s) \right\}, \quad \forall s \in \mathcal{S}, \mathbf{a} \in \mathbb{A},$$

924 from which it follows that
 925

$$V^\pi(s) = \min \left\{ \max \left\{ r_{\text{RAA}}(s), V^\pi(f(s, \pi(s))) \right\}, q(s) \right\}, \quad \forall s \in \mathcal{S}, \quad (15)$$

926 and
 927

$$v^*(s) = \min \left\{ \max \left\{ r_{\text{RAA}}(s), \max_{a \in \mathcal{A}} v^*(f(s, a)) \right\}, q(s) \right\}, \quad \forall s \in \mathcal{S}. \quad (16)$$

931 Since $x \in S_{n+1} \setminus S_n$, then $\sigma(x) = n + 1$ by definition of σ , so $\pi(x) \in$
 932 $\arg \max_{a \in \mathcal{A}} v_n^*(f(x, a))$ by (9). Thus
 933

$$v_n^*(f(x, \pi(x))) = \max_{a \in \mathcal{A}} v_n^*(f(x, a)). \quad (17)$$

936 But then
 937

$$v_n^*(f(x, \pi(x))) \geq v_n^*(f(x, w)) \geq \alpha_n \geq \alpha_{n+1} > -\infty,$$

938 where the second inequality comes from (10), the third comes from (7), and the final
 939 inequality comes from (3) ($\mathcal{S} \setminus S_n$ is non-empty because $x \in \mathcal{S} \setminus S_n$). Thus $f(x, \pi(x)) \in S_n$
 940 by (5). It then follows from (10) that

$$V^\pi(f(x, \pi(x))) = v^*(f(x, \pi(x))) = v_n^*(f(x, \pi(x))). \quad (18)$$

943 Now, observe that for all $s \in S_n$ and $s' \in \mathcal{S} \setminus S_n$,
 944

$$v^*(s) = v_n^*(s) \geq \alpha_n \geq v^*(s') \geq -\infty = v_n^*(s'), \quad (19)$$

946 where the first equality and inequality are from (10), the second inequality is from (11),
 947 and the final equality is from (5). Moreover, $f(x, a) \in S_n$ for at least one a (in particular
 948 $a = w$). Letting $\mathcal{A}' = \{a \in \mathcal{A} \mid f(x, a) \in S_n\}$, it follows from (19) that

$$\max_{a \in \mathcal{A}} v^*(f(x, a)) = \max_{a \in \mathcal{A}'} v^*(f(x, a)) = \max_{a \in \mathcal{A}'} v_n^*(f(x, a)) = \max_{a \in \mathcal{A}} v_n^*(f(x, a)). \quad (20)$$

951 From (17)-(20) we have
 952

$$V^\pi(f(x, \pi(x))) = \max_{a \in \mathcal{A}} v^*(f(x, a)) = \max_{a \in \mathcal{A}} v_n^*(f(x, a)). \quad (21)$$

955 Now observe that
 956

$$\begin{aligned} V^\pi(x) &= \min \left\{ \max \left\{ r_{\text{RAA}}(x), V^\pi(f(x, \pi(x))) \right\}, q(x) \right\}, \\ v^*(x) &= \min \left\{ \max \left\{ r_{\text{RAA}}(x), \max_{a \in \mathcal{A}} v^*(f(x, a)) \right\}, q(x) \right\}, \\ \alpha_{n+1} &= \min \left\{ \max \left\{ r_{\text{RAA}}(x), \max_{a \in \mathcal{A}} v_n^*(f(x, a)) \right\}, q(x) \right\}, \end{aligned}$$

963 where the first equation is from (15), the second is from (16), and the third is from (4). But
 964 then (14) follows from the above equations together with (21).

- 965 3. Let $x \in \mathcal{S} \setminus S_n$. We claim that $v^*(x) \leq \alpha_{n+1}$. Suppose otherwise. Then we can choose
 966 $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$ such that
 967

$$\min \left\{ r_{\text{RAA}}(\xi_x^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_x^{\mathbf{a}}(\kappa)) \right\} > \alpha_{n+1}. \quad (22)$$

970 It follows that $\xi_x^{\mathbf{a}}(\tau) \in S_n$, for otherwise
 971

$$\alpha_{n+1} \geq \min \{r_{\text{RAA}}(\xi_x^{\mathbf{a}}(\tau)), q(\xi_x^{\mathbf{a}}(\tau))\}$$

972 by (3), creating a contradiction.
 973
 974 So $x \notin S_n$ and $\xi_x^{\mathbf{a}}(\tau) \in S_n$, indicating that there is some $\theta \in \{0, \dots, \tau - 1\}$ such that
 975 $\xi_x^{\mathbf{a}}(\theta) \notin S_n$ and $f(\xi_x^{\mathbf{a}}(\theta), \mathbf{a}(\theta)) = \xi_x^{\mathbf{a}}(\theta + 1) \in S_n$. Moreover, $q(\xi_x^{\mathbf{a}}(\theta)) > \alpha_{n+1}$ by (22).
 976 It follows from claim 1 that $\xi_x^{\mathbf{a}}(\theta) \in S_{n+1}$.

977 But then it follows from claim 2 that $v^*(\xi_x^{\mathbf{a}}(\theta)) = \alpha_{n+1}$. However,
 978

$$\begin{aligned} v^*(\xi_x^{\mathbf{a}}(\theta)) &\geq \min \left\{ r_{\text{RAA}} \left(\xi_x^{\mathbf{a}}(\theta)(\tau - \theta) \right), \min_{\kappa \leq \tau - \theta} q \left(\xi_x^{\mathbf{a}}(\theta)(\kappa) \right) \right\} \\ &= \min \left\{ r_{\text{RAA}} (\xi_x^{\mathbf{a}}(\tau - \theta + \theta)), \min_{\kappa \leq \tau - \theta} q (\xi_x^{\mathbf{a}}(\kappa + \theta)) \right\} \\ &= \min \left\{ r_{\text{RAA}} (\xi_x^{\mathbf{a}}(\tau)), \min_{\kappa \in \{\theta, \theta+1, \dots, \tau\}} q (\xi_x^{\mathbf{a}}(\kappa)) \right\} \\ &> \alpha_{n+1}, \end{aligned}$$

987 giving the desired contradiction.
 988
 989

990 Having established these claims, we return to proving (12) and (13) hold. In fact, (13) follows
 991 immediately from claim 3, so we actually only need to show (12).

992 If $y \in S_n$, then from (5) and (10), we have that $V^\pi(y) = v^*(y) = v_n^*(y) = v_{n+1}^*(y)$, and from (7)
 993 and (10), we also have that $v_n^*(y) \geq \alpha_n \geq \alpha_{n+1}$. Together these establish (12) when $y \in S_n$.
 994

995 So suppose $y \in S_{n+1} \setminus S_n$. First, observe that $v_{n+1}^*(y) = \alpha_{n+1}$ by (5). There are now two
 996 possibilities. If there is some $a \in \mathcal{A}$ for which $f(y, a) \in S_n$, then (12) follows from claim 2. If
 997 instead, $f(y, a) \notin S_n$ for each $a \in \mathcal{A}$, then $\max_{a \in \mathcal{A}} v_n^*(f(y, a)) = -\infty$ by (5) (or if $n = 0$ by
 998 definition of v_0^*). Thus $\alpha_{n+1} = \min \{r_{\text{RAA}}(y), q(y)\}$ by (4), so

$$v^*(y) \geq V^\pi(y) \geq \min \{r_{\text{RAA}}(y), q(y)\} = \alpha_{n+1} \geq v^*(y),$$

1000 where the final inequality follows from claim 3. This completes the proof. \square
 1001
 1002

1003 **Corollary 4.** For all $s \in \mathcal{S}$, we have $\tilde{V}_{\text{RA}}^*(s) = \tilde{v}_{\text{RA}}^*(s)$.
 1004

1005 **Lemma 3.** Let $F : \mathbb{A} \times \mathbb{N} \rightarrow \mathbb{R}$. Then

$$\sup_{\mathbf{a} \in \mathbb{A}} \sup_{\tau \in \mathbb{N}} \sup_{\mathbf{a}' \in \mathbb{A}'} F([\mathbf{a}, \mathbf{a}']_\tau, \tau) = \sup_{\mathbf{a} \in \mathbb{A}} \sup_{\tau \in \mathbb{N}} F(\mathbf{a}, \tau). \quad (23)$$

1009 *Proof.* We proceed by showing both inequalities corresponding to (23) hold.
 1010

1011
 1012 (\geq) Given any $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$, we have $\sup_{\mathbf{a}' \in \mathbb{A}'} F([\mathbf{a}, \mathbf{a}']_\tau, \tau) \geq F(\mathbf{a}, \tau)$. Taking the
 1013 suprema over $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$ on both sides of this inequality gives the desired result.
 1014

1015 (\leq) Given any $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$, we have

$$\sup_{\mathbf{a}' \in \mathbb{A}'} F([\mathbf{a}, \mathbf{a}']_\tau, \tau) \leq \sup_{\mathbf{a}'' \in \mathbb{A}} F(\mathbf{a}'', \tau),$$

1019 so that the result follows from taking the suprema over $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$ on both sides of
 1020 this inequality.
 1021
 1022 \square
 1023

1024 **Lemma 4.** For each $s \in \mathcal{S}$,

$$v_{\text{RAA}}^*(s) = \tilde{v}_{\text{RA}}^*(s).$$

1026 *Proof.* For each $s \in \mathcal{S}$, we have
1027

$$\tilde{v}_{\text{RA}}^*(s) = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r_{\text{RAA}}(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} \quad (24)$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ r(\xi_s^{\mathbf{a}}(\tau)), v_{\text{A}}^*(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} \quad (25)$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \min_{\kappa' \in \mathbb{N}} q\left(\xi_{\xi_s^{\mathbf{a}}(\tau)}^{\mathbf{a}'}(\kappa')\right), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \min_{\kappa' \in \mathbb{N}} q\left(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \kappa')\right), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa' \in \mathbb{N}} q\left(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \kappa')\right), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \min \left\{ r\left(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau)\right), \min_{\kappa' \in \mathbb{N}} q\left(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \kappa')\right), \min_{\kappa \leq \tau} q\left(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\kappa)\right) \right\} \quad (26)$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa' \in \mathbb{N}} q(\xi_s^{\mathbf{a}}(\tau + \kappa')), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} \quad (27)$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$= \max_{\mathbf{a} \in \mathbb{A}} \left\{ \max_{\tau \in \mathbb{N}} r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$= v_{\text{RAA}}^*(s),$$

1051 where the equality between (24) and (25) follows from Corollary 3, and where the equality between
1052 (26) and (27) follows from Lemma 3. \square
1053

1054 Before the next lemma, we need to introduce two last pieces of notation. First, we let $\bar{\Pi}$ be the set of
1055 augmented policies $\bar{\pi} : \mathcal{S} \times \mathcal{Y} \times \mathcal{Z} \rightarrow \mathcal{A}$, where
1056

$$\mathcal{Y} = \{r(s) \mid s \in \mathcal{S}\} \quad \text{and} \quad \mathcal{Z} = \{q(s) \mid s \in \mathcal{S}\}.$$

1058 Next, given $s \in \mathcal{S}$, $y \in \mathcal{Y}$, $z \in \mathcal{Z}$, and $\bar{\pi} \in \bar{\Pi}$, we let $\bar{\xi}_s^{\bar{\pi}} : \mathbb{N} \rightarrow \mathcal{S}$, $\bar{\eta}_s^{\bar{\pi}} : \mathbb{N} \rightarrow \mathcal{Y}$, and $\bar{\zeta}_s^{\bar{\pi}} : \mathbb{N} \rightarrow \mathcal{Z}$, be
1059 the solution of the evolution
1060

$$\begin{aligned} \bar{\xi}_s^{\bar{\pi}}(t+1) &= f\left(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\pi}\left(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\eta}_s^{\bar{\pi}}(t), \bar{\zeta}_s^{\bar{\pi}}(t)\right)\right), \\ \bar{\eta}_s^{\bar{\pi}}(t+1) &= \max\left\{r\left(\bar{\xi}_s^{\bar{\pi}}(t+1)\right), \bar{\eta}_s^{\bar{\pi}}(t)\right\}, \\ \bar{\zeta}_s^{\bar{\pi}}(t+1) &= \min\left\{q\left(\bar{\xi}_s^{\bar{\pi}}(t+1)\right), \bar{\zeta}_s^{\bar{\pi}}(t)\right\}, \end{aligned}$$

1065 for which $\bar{\xi}_s^{\bar{\pi}}(0) = s$, $\bar{\eta}_s^{\bar{\pi}}(0) = r(s)$, and $\bar{\zeta}_s^{\bar{\pi}}(0) = q(s)$.
1066

1067 **Lemma 5.** *There is a $\bar{\pi} \in \bar{\Pi}$ such that*

$$v_{\text{RAA}}^*(s) = \min \left\{ \max_{\tau \in \mathbb{N}} r\left(\bar{\xi}_s^{\bar{\pi}}(\tau)\right), \min_{\tau \in \mathbb{N}} q\left(\bar{\xi}_s^{\bar{\pi}}(\tau)\right) \right\} \quad (28)$$

1070 for all $s \in \mathcal{S}$.
1071

1072 *Proof.* By Lemmas 1 and 2 together with Corollary 3, we can choose $\pi, \theta \in \Pi$ such that
1073

$$\tilde{v}_{\text{RA}}^*(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\pi}(\tau)), v_{\text{A}}^*(\xi_s^{\pi}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\pi}(\kappa)) \right\} \quad \forall s \in \mathcal{S},$$

$$v_{\text{A}}^*(s) = \min_{\tau \in \mathbb{N}} q(\xi_s^{\theta}(\tau)) \quad \forall s \in \mathcal{S}.$$

1078 We introduce some useful notation we will use throughout the rest of the proof. For each $s \in \mathcal{S}$, let
1079 $[s]^+ = f(s, \pi(s))$, $[y]_s^+ = \max\{y, r([s]^+)\}$, $[z]_s^+ = \min\{z, q([s]^+)\}$.

1080 We define an augmented policy $\bar{\pi} \in \overline{\Pi}$ by
 1081

$$1082 \bar{\pi}(s, y, z) = \begin{cases} \pi(s) & \min\{[y]_s^+, [z]_s^+, v_A^*([s]^+)\} \geq \min\{y, z, v_A^*(s)\}, \\ 1083 \theta(s) & \text{otherwise.} \end{cases}$$

1085 Now fix some $s \in \mathcal{S}$. For all $t \in \mathbb{N}$, set $\bar{x}_t = \xi_s^{\bar{\pi}}(t)$, $\bar{y}_t = \bar{\eta}_s^{\bar{\pi}}(t) = \max_{\tau \leq t} r(\bar{x}_\tau)$, and $\bar{z}_t = \zeta_s^{\bar{\pi}}(t) = 1086 \min_{\tau \leq t} q(\bar{x}_\tau)$, and also set $x_t^\circ = \xi_s^\pi(t)$, $y_t^\circ = \max_{\tau \leq t} r(x_\tau^\circ)$, and $z_t^\circ = \min_{\tau \leq t} q(x_\tau^\circ)$.
 1087

1088 First, assume that t is such that $\min\{[\bar{y}_t]_{\bar{x}_t}^+, [\bar{z}_t]_{\bar{x}_t}^+, v_A^*([\bar{x}_t]^+)\} < \min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\}$. In this case,
 1089 $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \theta(\bar{x}_t)$, so that

$$1090 \min\{\bar{z}_t, v_A^*(\bar{x}_t)\} = \min\{\bar{z}_{t+1}, v_A^*(\bar{x}_{t+1})\}$$

1091 by our choice of θ . Since \bar{y}_t is non-decreasing in t , thus have
 1092

$$1093 \min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\} \leq \min\{\bar{y}_{t+1}, \bar{z}_{t+1}, v_A^*(\bar{x}_{t+1})\}.$$

1094 Next, assume that t is such that $\min\{[\bar{y}_t]_{\bar{x}_t}^+, [\bar{z}_t]_{\bar{x}_t}^+, v_A^*([\bar{x}_t]^+)\} \geq \min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\}$. In this case,
 1095 we have that $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \pi(\bar{x}_t)$, so
 1096

$$1097 \min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\} \leq \min\{[\bar{y}_t]_{\bar{x}_t}^+, [\bar{z}_t]_{\bar{x}_t}^+, v_A^*([\bar{x}_t]^+)\} = \min\{\bar{y}_{t+1}, \bar{z}_{t+1}, v_A^*(\bar{x}_{t+1})\}.$$

1098 It thus follows from these two cases that $\min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\}$ is non-decreasing in t . Let
 1099

$$1100 T = \min \{t \in \mathbb{N} \mid \min\{[\bar{y}_t]_{\bar{x}_t}^+, [\bar{z}_t]_{\bar{x}_t}^+, v_A^*([\bar{x}_t]^+)\} < \min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\}\}.$$

1101 There are again two cases:
 1102

1103 (1) $(T < \infty)$ In this case, $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \pi(\bar{x}_t)$ for $t < T$. Then $\bar{x}_t = x_t^\circ$, $\bar{y}_t = y_t^\circ$, and $\bar{z}_t = z_t^\circ$ for all
 1104 $t \leq T$. It follows that $[\bar{x}_t]^+ = x_{t+1}^\circ$, $[\bar{y}_t]_{\bar{x}_t}^+ = y_{t+1}^\circ$, and $[\bar{z}_t]_{\bar{x}_t}^+ = z_{t+1}^\circ$ for all $t \leq T$. Thus
 1105 by definition of T ,

$$1106 \min\{y_{t+1}^\circ, z_{t+1}^\circ, v_A^*(x_{t+1}^\circ)\} \geq \min\{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\} \quad \forall t < T.$$

1107 and
 1108

$$\min\{y_{T+1}^\circ, z_{T+1}^\circ, v_A^*(x_{T+1}^\circ)\} < \min\{y_T^\circ, z_T^\circ, v_A^*(x_T^\circ)\}.$$

1109 But since y_t° is non-decreasing and $\min\{z_t^\circ, v_A^*(x_t^\circ)\}$ is non-increasing in t , it follows that
 1110 $\min\{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\}$ must achieve its maximal value at the smallest t for which it strictly
 1111 decreases from t to $t+1$, i.e.
 1112

$$1113 \begin{aligned} \min\{\bar{y}_T, \bar{z}_T, v_A^*(\bar{x}_T)\} &= \min\{y_T^\circ, z_T^\circ, v_A^*(x_T^\circ)\} \\ 1114 &= \max_{t \in \mathbb{N}} \min\{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\} \\ 1115 &\geq \max_{t \in \mathbb{N}} \min\{r(x_t^\circ), z_t^\circ, v_A^*(x_t^\circ)\} \\ 1116 &= \tilde{v}_{RA}^*(s). \end{aligned}$$

1117 where the final equality follows from our choice of π . Since $\min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\}$ is non-decreasing in t , then
 1118

$$1119 \min\{\bar{y}_t, \bar{z}_t\} \geq \min\{\bar{y}_t, \bar{z}_t, v_A^*(\bar{x}_t)\} \geq \min\{\bar{y}_T, \bar{z}_T, v_A^*(\bar{x}_T)\} = \tilde{v}_{RA}^*(s) \quad \forall t \geq T.$$

1120 Thus
 1121

$$1122 v_{RAA}^*(s) \geq \min \left\{ \max_{t \in \mathbb{N}} r(\bar{x}_t), \min_{t \in \mathbb{N}} q(\bar{x}_t) \right\} = \lim_{t \rightarrow \infty} \min\{\bar{y}_t, \bar{z}_t\} \geq \tilde{v}_{RA}^*(s) = v_{RAA}^*(s),$$

1123 where the final equality follows from Lemma (4). Thus the proof is complete in this case.
 1124

1125 (2) $(T = \infty)$ In this case, $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \pi(\bar{x}_t)$ for all $t \in \mathbb{N}$. Then $\bar{x}_t = x_t^\circ$, $\bar{y}_t = y_t^\circ$, and $\bar{z}_t = z_t^\circ$ for
 1126 all $t \in \mathbb{N}$. Also $[\bar{x}_t]^+ = x_{t+1}^\circ$, $[\bar{y}_t]_{\bar{x}_t}^+ = y_{t+1}^\circ$, and $[\bar{z}_t]_{\bar{x}_t}^+ = z_{t+1}^\circ$ for all $t \in \mathbb{N}$. Thus by
 1127 definition of T ,

$$1128 \min\{y_{t+1}^\circ, z_{t+1}^\circ, v_A^*(x_{t+1}^\circ)\} \geq \min\{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\} \quad \forall t \in \mathbb{N}.$$

1134 Let $T' \in \arg \max_{t \in \mathbb{N}} \min \{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\}$. Then
 1135

$$\begin{aligned} 1136 \quad \min \{\bar{y}_{T'}, \bar{z}_{T'}, v_A^*(\bar{x}_{T'})\} &= \min \{y_{T'}^\circ, z_{T'}^\circ, v_A^*(x_{T'}^\circ)\} \\ 1137 \quad &= \max_{t \in \mathbb{N}} \min \{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\} \\ 1138 \quad &\geq \max_{t \in \mathbb{N}} \min \{r(x_t^\circ), z_t^\circ, v_A^*(x_t^\circ)\} \\ 1139 \quad &= \tilde{v}_{RA}^*(s). \end{aligned}$$

1140 The rest of the proof follows the same as the previous case with T replaced by T' .
 1141

□

1142 **Corollary 5.** For all $s \in \mathcal{S}$, we have $V_{RAA}^*(s) = v_{RAA}^*(s)$.
 1143

1144 *Proof of Theorem 1.* Theorem 1 is now a direct consequence of the previous corollary together with
 1145 Corollary 4 and Lemma 4. □
 1146

1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

1188 **B PROOF OF RR MAIN THEOREM**
 1189

1190 We first define the value functions, $V_{R1}^*, V_{R2}^*, \tilde{V}_R^*, V_{RR}^* : \mathcal{S} \rightarrow \mathbb{R}$ by
 1191

$$\begin{aligned} V_{R1}^*(s) &= \max_{\pi \in \Pi} \max_{\tau \in \mathbb{N}} r_1(\xi_s^\pi(\tau)), \\ V_{R2}^*(s) &= \max_{\pi \in \Pi} \max_{\tau \in \mathbb{N}} r_2(\xi_s^\pi(\tau)), \\ \tilde{V}_R^*(s) &= \max_{\pi \in \Pi} \max_{\tau \in \mathbb{N}} r_{RR}(\xi_s^\pi(\tau)), \\ V_{RR}^*(s) &= \max_{\pi \in \Pi} \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\xi_s^\pi(\tau)), \max_{\tau \in \mathbb{N}} r_2(\xi_s^\pi(\tau)) \right\}. \end{aligned}$$

1200 We next define the value functions, $v_{R1}^*, v_{R2}^*, \tilde{v}_R^*, v_{RR}^* : \mathcal{S} \rightarrow \mathbb{R}$, which maximize over action sequences
 1201 rather than policies:
 1202

$$\begin{aligned} v_{R1}^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\mathbf{a}}(\tau)), \\ v_{R2}^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\mathbf{a}}(\tau)), \\ \tilde{v}_R^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} r_{RR}(\xi_s^{\mathbf{a}}(\tau)), \\ v_{RR}^*(s) &= \max_{\mathbf{a} \in \mathbb{A}} \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\mathbf{a}}(\tau)) \right\}, \end{aligned}$$

1210 where r_{RR} is as in Theorem 2. Observe that for each $s \in \mathcal{S}$,
 1211

$$v_{R1}^*(s) \geq V_{R1}^*(s), \quad v_{R2}^*(s) \geq V_{R2}^*(s), \quad \tilde{v}_R^*(s) \geq \tilde{V}_R^*(s), \quad v_{RR}^*(s) \geq V_{RR}^*(s).$$

1214 We now prove a series of lemmas that will be useful in the proof of the main theorem.
 1215

Lemma 6. *There are $\pi_1, \pi_2 \in \Pi$ such that*

$$v_{R1}^*(s) = \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\pi_1}(\tau)) \text{ and } v_{R2}^*(s) = \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\pi_2}(\tau))$$

1219 for all $s \in \mathcal{S}$.
 1220

1221 *Proof.* We will just prove the result for $v_{R1}^*(s)$ since the other result follows identically. For each
 1222 $s \in \mathcal{S}$, let τ_s be the smallest element of \mathbb{N} for which
 1223

$$\max_{\mathbf{a} \in \mathbb{A}} r_1(\xi_s^{\mathbf{a}}(\tau_s)) = v_{R1}^*(s).$$

1226 Moreover, for each $s \in \mathcal{S}$, let \mathbf{a}_s be such that
 1227

$$r_1(\xi_s^{\mathbf{a}_s}(\tau_s)) = v_{R1}^*(s).$$

1229 Let $\pi_1 \in \Pi$ be given by $\pi_1(s) = \mathbf{a}_s(0)$. It suffices to show that
 1230

$$r_1(\xi_s^{\pi_1}(\tau_s)) = v_{R1}^*(s) \tag{29}$$

1232 for all $s \in \mathcal{S}$, for in this case, we have
 1233

$$v_{R1}^*(s) \geq \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\pi_1}(\tau)) \geq r_1(\xi_s^{\pi_1}(\tau_s)) = v_{R1}^*(s) \quad \forall s \in \mathcal{S}.$$

1236 We show (29) holds for each $s \in \mathcal{S}$ by induction on τ_s . First, suppose that $s \in \mathcal{S}$ is such that $\tau_s = 0$.
 1237 Then

$$r_1(\xi_s^{\pi_1}(\tau_s)) = r_1(s) = r_1(\xi_s^{\mathbf{a}_s}(\tau_s)) = v_{R1}^*(s).$$

1240 For the induction step, let $n \in \mathbb{N}$ and suppose that
 1241

$$r_1(\xi_s^{\pi_1}(\tau_s)) = v_{R1}^*(s) \quad \forall s \in \mathcal{S} \text{ such that } \tau_s \leq n.$$

1242 Now fix some $x \in \mathcal{S}$ such that $\tau_x = n + 1$. Notice that
 1243

$$\begin{aligned} v_{\text{R1}}^*(x) &\geq v_{\text{R1}}^*(f(x, \pi_1(x))) \\ &\geq \max_{\mathbf{a} \in \mathbb{A}} r_1(\xi_{f(x, \pi_1(x))}^{\mathbf{a}}(n)) \\ &\geq r_1(\xi_{f(x, \pi_1(x))}^{\mathbf{a}_x|_1}(n)) \\ &= r_1(\xi_x^{[\pi_1(x), \mathbf{a}_x|_1]}(n+1)) \\ &= r_1(\xi_x^{\mathbf{a}_x}(\tau_x)) \\ &= v_{\text{R1}}^*(x), \end{aligned}$$

1253 so that $v_{\text{R1}}^*(f(x, \pi_1(x))) = v_{\text{R1}}^*(x)$ and $\tau_{f(x, \pi_1(x))} \leq n$. It suffices to show
 1254

$$\tau_{f(x, \pi_1(x))} = n, \quad (30)$$

1255 for then, by the induction assumption, we have
 1256

$$r_1(\xi_x^{\pi_1}(\tau_x)) = r_1(\xi_{f(x, \pi_1(x))}^{\pi_1}(n)) = v_{\text{R1}}^*(f(x, \pi_1(x))) = v_{\text{R1}}^*(x).$$

1260 To show (30), assume instead that
 1261

$$\tau_{f(x, \pi_1(x))} < n.$$

1262 But
 1263

$$\begin{aligned} v_{\text{R1}}^*(x) &\geq \max_{\mathbf{a} \in \mathbb{A}} r_1(\xi_x^{\mathbf{a}}(\tau_{f(x, \pi_1(x))} + 1)) \\ &\geq r_1(\xi_x^{[\pi_1(x), \mathbf{a}_{f(x, \pi_1(x))}]}(\tau_{f(x, \pi_1(x))} + 1)) \\ &= r_1(\xi_{f(x, \pi_1(x))}^{\mathbf{a}_{f(x, \pi_1(x))}}(\tau_{f(x, \pi_1(x))})) \\ &= v_{\text{R1}}^*(f(x, \pi_1(x))) \\ &= v_{\text{R1}}^*(x), \end{aligned}$$

1272 so that
 1273

$$v_{\text{R1}}^*(x) = \max_{\mathbf{a} \in \mathbb{A}} r_1(\xi_x^{\mathbf{a}}(\tau_{f(x, \pi_1(x))} + 1))$$

1275 and thus
 1276

$$\tau_x \leq \tau_{f(x, \pi_1(x))} + 1 < n + 1,$$

1277 giving our desired contradiction. \square

1278 **Corollary 6.** For all $s \in \mathcal{S}$, we have $V_{\text{R1}}^*(s) = v_{\text{R1}}^*(s)$ and $V_{\text{R2}}^*(s) = v_{\text{R2}}^*(s)$.

1279 **Lemma 7.** There is a $\pi \in \Pi$ such that
 1280

$$\tilde{v}_{\text{R}}^*(s) = \max_{\tau \in \mathbb{N}} r_{\text{RR}}(\xi_s^{\pi}(\tau)).$$

1283 for all $s \in \mathcal{S}$.
 1284

1285 *Proof.* This lemma follows by precisely the same proof as the previous lemma, with r_1 , v_{R1}^* , and π_1
 1286 replaced with r_{RR} , \tilde{v}_{R}^* , and π respectively. \square

1287 **Corollary 7.** For all $s \in \mathcal{S}$, we have $\tilde{V}_{\text{R}}^*(s) = \tilde{v}_{\text{R}}^*(s)$.

1288 **Lemma 8.** Let $\zeta_1 : \mathbb{N} \rightarrow \mathbb{R}$ and $\zeta_2 : \mathbb{N} \rightarrow \mathbb{R}$. Then
 1289

$$\begin{aligned} &\sup_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}, \min \left\{ \sup_{\tau' \in \mathbb{N}} \zeta_1(\tau + \tau'), \zeta_2(\tau) \right\} \right\} \\ &= \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau), \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) \right\}. \end{aligned}$$

1295 *Proof.* We proceed by showing both inequalities corresponding to the above equality hold.

1296 (\leq) Observe that
 1297
 1298
 1299
 1300 $\sup_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}, \min \left\{ \sup_{\tau' \in \mathbb{N}} \zeta_1(\tau + \tau'), \zeta_2(\tau) \right\} \right\}$
 1301
 1302 $\leq \max \left\{ \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau), \sup_{\tau \in \mathbb{N}} \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}, \min \left\{ \sup_{\tau \in \mathbb{N}} \sup_{\tau' \in \mathbb{N}} \zeta_1(\tau + \tau'), \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) \right\} \right\}$
 1303
 1304 $= \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau), \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) \right\}$
 1305
 1306
 1307
 1308
 1309
 1310 (\geq) Fix $\varepsilon > 0$. Choose $\tau_1, \tau_2 \in \mathbb{N}$ such that $\zeta_1(\tau_1) \geq \sup_{\tau \in \mathbb{N}} \zeta_1(\tau) - \varepsilon$ and $\zeta_2(\tau_2) \geq \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) - \varepsilon$. Without loss of generality, we can assume $\tau_1 \leq \tau_2$. Then
 1311
 1312
 1313
 1314
 1315 $\sup_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}, \min \left\{ \sup_{\tau' \in \mathbb{N}} \zeta_1(\tau + \tau'), \zeta_2(\tau) \right\} \right\}$
 1316
 1317 $\geq \sup_{\tau \in \mathbb{N}} \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}$
 1318
 1319 $\geq \min \left\{ \zeta_1(\tau_1), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau_1 + \tau') \right\}$
 1320
 1321 $\geq \min \{ \zeta_1(\tau_1), \zeta_2(\tau_2) \}$
 1322
 1323 $\geq \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau) - \varepsilon, \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) - \varepsilon \right\}$
 1324
 1325 $= \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau), \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) \right\} - \varepsilon.$
 1326
 1327
 1328
 1329
 1330 But since $\varepsilon > 0$ was arbitrary, the desired inequality follows.
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338 \square
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346 **Lemma 9.** For each $s \in \mathcal{S}$,
 1347
 1348
 1349 $\tilde{v}_R^*(s) = v_{RR}^*(s).$

1350 *Proof.* For each $s \in \mathcal{S}$,

$$1352 \quad \tilde{v}_R^*(s) = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} r_{RR}(\xi_s^{\mathbf{a}}(\tau)) \quad (31)$$

$$1353 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \{ \min \{ r_1(\xi_s^{\mathbf{a}}(\tau)), v_{R2}^*(\xi_s^{\mathbf{a}}(\tau)) \}, \min \{ v_{R1}^*(\xi_s^{\mathbf{a}}(\tau)), r_2(\xi_s^{\mathbf{a}}(\tau)) \} \} \quad (32)$$

$$1354 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_2(\xi_{\xi_s^{\mathbf{a}}(\tau)}^{\mathbf{a}'}(\tau')) \right\}, \right.$$

$$1355 \quad \left. \min \left\{ \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_1(\xi_{\xi_s^{\mathbf{a}}(\tau)}^{\mathbf{a}'}(\tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\}$$

$$1356 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')) \right\}, \right.$$

$$1357 \quad \left. \min \left\{ \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\}$$

$$1358 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')) \right\}, \right.$$

$$1359 \quad \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\}$$

$$1360 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \max_{\mathbf{a}'' \in \mathbb{A}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')) \right\}, \right.$$

$$1361 \quad \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\}$$

$$1362 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \max_{\mathbf{a}'' \in \mathbb{A}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')) \right\}, \right.$$

$$1363 \quad \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \quad (33)$$

$$1364 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')) \right\}, \right.$$

$$1365 \quad \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\}$$

$$1366 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \max_{\mathbf{a}'' \in \mathbb{A}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')) \right\}, \right.$$

$$1367 \quad \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \quad (34)$$

$$1368 \quad = \max_{\mathbf{a} \in \mathbb{A}} \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \quad (35)$$

$$1369 \quad = v_{RR}^*(s),$$

1370 where the equality between 31 and 32 follows from Corollary 6, the equality between 33 and 34
1371 follows from Lemma 3, and the equality between 34 and 35 follows from Lemma 8. \square

1372 Before the next lemma, we need to introduce two last pieces of notation. First, we let $\bar{\Pi}$ be the set of
1373 augmented policies $\bar{\pi} : \mathcal{S} \times \mathcal{Y} \times \mathcal{Z} \rightarrow \mathcal{A}$, as in the previous section, but where
1374

$$1375 \quad \mathcal{Y} = \{r_1(s) \mid s \in \mathcal{S}\} \quad \text{and} \quad \mathcal{Z} = \{r_2(s) \mid s \in \mathcal{S}\}.$$

1376 Next, given $s \in \mathcal{S}$, $y \in \mathcal{Y}$, $z \in \mathcal{Z}$, and $\bar{\pi} \in \bar{\Pi}$, we let $\bar{\xi}_s^{\bar{\pi}} : \mathbb{N} \rightarrow \mathcal{S}$, $\bar{\eta}_s^{\bar{\pi}} : \mathbb{N} \rightarrow \mathcal{Y}$, and $\bar{\zeta}_s^{\bar{\pi}} : \mathbb{N} \rightarrow \mathcal{Z}$, be
1377 the solution of the evolution

$$1378 \quad \bar{\xi}_s^{\bar{\pi}}(t+1) = f(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\pi}(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\eta}_s^{\bar{\pi}}(t), \bar{\zeta}_s^{\bar{\pi}}(t))),$$

$$1379 \quad \bar{\eta}_s^{\bar{\pi}}(t+1) = \max \{r_1(\bar{\xi}_s^{\bar{\pi}}(t+1)), \bar{\eta}_s^{\bar{\pi}}(t)\},$$

$$1380 \quad \bar{\zeta}_s^{\bar{\pi}}(t+1) = \max \{r_2(\bar{\xi}_s^{\bar{\pi}}(t+1)), \bar{\zeta}_s^{\bar{\pi}}(t)\},$$

1381 for which $\bar{\xi}_s^{\bar{\pi}}(0) = s$, $\bar{\eta}_s^{\bar{\pi}}(0) = r_1(s)$, and $\bar{\zeta}_s^{\bar{\pi}}(0) = r_2(s)$.

1382 **Lemma 10.** *There is a $\bar{\pi} \in \bar{\Pi}$ such that*

$$1383 \quad v_{RR}^*(s) = \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\bar{\xi}_s^{\bar{\pi}}(\tau)), \max_{\tau \in \mathbb{N}} r_2(\bar{\xi}_s^{\bar{\pi}}(\tau)) \right\}$$

1384 for all $s \in \mathcal{S}$.

1404 *Proof.* By Lemmas 6 and 7 together with Corollary 6, we can choose $\pi, \theta_1, \theta_2 \in \Pi$ such that
 1405

$$v_{R1}^*(s) = \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\theta_1}(\tau)) \quad \forall s \in \mathcal{S},$$

$$v_{R2}^*(s) = \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\theta_2}(\tau)) \quad \forall s \in \mathcal{S},$$

$$\tilde{v}_R^*(s) = \max_{\tau \in \mathbb{N}} \max \{ \min \{ r_1(\xi_s^\pi(\tau)), v_{R2}^*(\xi_s^\pi(\tau)) \}, \min \{ r_2(\xi_s^\pi(\tau)), v_{R1}^*(\xi_s^\pi(\tau)) \} \} \quad \forall s \in \mathcal{S}.$$

1411 Define $\bar{\pi} \in \overline{\Pi}$ by

$$\bar{\pi}(s, y, z) = \begin{cases} \pi(s) & \max\{y, z\} < \tilde{v}_R^*(s) \\ \theta_1(s) & \max\{y, z\} \geq \tilde{v}_R^*(s) \text{ and } y \leq z, \\ \theta_2(s) & \max\{y, z\} \geq \tilde{v}_R^*(s) \text{ and } y > z. \end{cases}$$

1416 Now fix some $s \in \mathcal{S}$. For all $t \in \mathbb{N}$, set $\bar{x}_t = \bar{\xi}_s^{\bar{\pi}}(t)$, $\bar{y}_t = \bar{\eta}_s^{\bar{\pi}}(t) = \max_{\tau \leq t} r_1(\bar{x}_\tau)$, and $\bar{z}_t = \bar{\zeta}_s^{\bar{\pi}}(t) = \max_{\tau \leq t} r_2(\bar{x}_\tau)$, and also set $x_t^\circ = \xi_s^\pi(t)$. It suffices to show
 1417

$$v_{RR}^*(s) \leq \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\bar{x}_\tau), \max_{\tau \in \mathbb{N}} r_2(\bar{x}_\tau) \right\}, \quad (36)$$

1421 since the reverse inequality is immediate. We proceed in three steps.
 1422

- 1423 1. We claim there exists a $t \in \mathbb{N}$ such that $\max \{r_1(\bar{x}_t), r_2(\bar{x}_t)\} \geq \tilde{v}_R^*(\bar{x}_t)$.

1424 Suppose otherwise. Then $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \pi(\bar{x}_t)$ so that $\bar{x}_t = x_t^\circ$ for all $t \in \mathbb{N}$. Thus

$$\begin{aligned} \max_{t \in \mathbb{N}} \max \{r_1(\bar{x}_t), r_2(\bar{x}_t)\} &< \max_{t \in \mathbb{N}} \tilde{v}_R^*(\bar{x}_t) \\ &= \tilde{v}_R^*(s) \\ &= \max_{\tau \in \mathbb{N}} \max \{ \min \{r_1(x_\tau^\circ), v_{R2}^*(x_\tau^\circ)\}, \min \{r_2(x_\tau^\circ), v_{R1}^*(x_\tau^\circ)\} \} \\ &= \max_{\tau \in \mathbb{N}} \max \{ \min \{r_1(\bar{x}_\tau), v_{R2}^*(\bar{x}_\tau)\}, \min \{r_2(\bar{x}_\tau), v_{R1}^*(\bar{x}_\tau)\} \} \\ &\leq \max_{\tau \in \mathbb{N}} \max \{r_1(\bar{x}_\tau), r_2(\bar{x}_\tau)\}, \end{aligned}$$

1434 providing the desired contradiction.
 1435

- 1436 2. Let T be the smallest element of \mathbb{N} for which

$$\max \{r_1(\bar{x}_T), r_2(\bar{x}_T)\} \geq \tilde{v}_R^*(\bar{x}_T),$$

1439 which must exist by the previous step, and let T' be the smallest element of \mathbb{N} for which

$$\max \{ \min \{r_1(x_{T'}^\circ), v_{R2}^*(x_{T'}^\circ)\}, \min \{r_2(x_{T'}^\circ), v_{R1}^*(x_{T'}^\circ)\} \} = \tilde{v}_R^*(s),$$

1442 which must exist by our choice of π . We claim $T' \geq T$.

1443 Suppose otherwise. Since $\bar{x}_t = x_t^\circ$ for all $t \leq T$, then in particular $\bar{x}_{T'} = x_{T'}^\circ$, so that

$$\max \{ \min \{r_1(\bar{x}_{T'}), v_{R2}^*(\bar{x}_{T'})\}, \min \{r_2(\bar{x}_{T'}), v_{R1}^*(\bar{x}_{T'})\} \} = \tilde{v}_R^*(s).$$

1446 But then

$$\max \{r_1(\bar{x}_{T'}), r_2(\bar{x}_{T'})\} \geq \tilde{v}_R^*(s) \geq \tilde{v}_R^*(\bar{x}_{T'}).$$

1449 By our choice of T , we then have $T \leq T'$, creating a contradiction.
 1450

- 1451 3. It follows from the previous step that

$$\tilde{v}_R^*(\bar{x}_T) = \tilde{v}_R^*(x_T^\circ) = \tilde{v}_R^*(s).$$

1453 By our choice of T , there are two cases: $r_1(\bar{x}_T) \geq \tilde{v}_R^*(\bar{x}_T)$ and $r_2(\bar{x}_T) \geq \tilde{v}_R^*(\bar{x}_T)$. We
 1455 assume the first case and prove the desired result, with case two following identically. To
 1456 reach a contradiction, assume
 1457

$$r_2(\bar{x}_t) < \tilde{v}_R^*(\bar{x}_T) \quad \forall t \in \mathbb{N}.$$

1458 But then $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \theta_2(\bar{x}_t)$ for all $t \geq T$, so $v_{\text{R}}^*(\bar{x}_T) = \max_{t \geq T} r_2(\bar{x}_t) < \tilde{v}_{\text{R}}^*(\bar{x}_T) \leq$
 1459 $\tilde{v}_{\text{R}}^*(s)$. Thus $r_2(x_{T'}^\circ) \leq v_{\text{R}}^*(x_{T'}^\circ) \leq v_{\text{R}}^*(x_T^\circ) = v_{\text{R}}^*(\bar{x}_T) < \tilde{v}_{\text{R}}^*(s)$. It follows that
 1460

$$1461 \max \{ \min \{ r_1(x_{T'}^\circ), v_{\text{R}}^*(x_{T'}^\circ) \}, \min \{ r_2(x_{T'}^\circ), v_{\text{R}}^*(x_{T'}^\circ) \} \} < \tilde{v}_{\text{R}}^*(s),$$

1462 contradicting our choice of T' .

1463 Thus $r_2(\bar{x}_t) \geq \tilde{v}_{\text{R}}^*(\bar{x}_T) = \tilde{v}_{\text{R}}^*(s)$ for some $t \in \mathbb{N}$ and also $r_1(\bar{x}_T) \geq \tilde{v}_{\text{R}}^*(\bar{x}_T) = \tilde{v}_{\text{R}}^*(s)$, so
 1464 that (36) must hold by Lemma 9.

□

1466
 1467 **Corollary 8.** For all $s \in \mathcal{S}$, we have $V_{\text{RR}}^*(s, r_1(s), r_2(s)) = v_{\text{RR}}^*(s)$.

1468
 1469 *Proof of Theorem 2.* The proof of this theorem immediately follows from the previous corollary
 1470 together with Corollary 7 and Lemma 9. □
 1471

C PROOF OF OPTIMALITY THEOREM

1475 *Proof of Theorem 3.* The inequalities in both lines of the theorem follow from the fact that for each
 1476 $\pi \in \Pi$, we can define a corresponding augmented policy $\bar{\pi} \in \bar{\Pi}$ by

$$1477 \bar{\pi}(s, y, z) = \pi(s) \quad \forall s \in \mathcal{S}, y \in \mathcal{Y}, z \in \mathcal{Z},$$

1478 in which case $V_{\text{RAA}}^\pi(s) = V_{\text{RAA}}^{\bar{\pi}}(s)$ and $V_{\text{RR}}^\pi(s) = V_{\text{RR}}^{\bar{\pi}}(s)$ for each $s \in \mathcal{S}$. Note that in general, we
 1479 cannot define a corresponding policy for each augmented policy, so the reverse inequality does not
 1480 generally hold (see Figure 3 for intuition regarding this fact).

1481 The equalities in both lines of the theorem are simply restatements of Lemma 5 and Lemma 9. □
 1482

D THE SRABE AND ITS POLICY GRADIENT

1486 *Proof of Proposition 1.* We here closely follow the proof of Theorem 3 in So et al. (2024), which
 1487 itself modifies the proofs of the Policy Gradient Theorems in Chapter 13.2 and 13.6 Sutton and Barto
 1488 (2018). We only make the minimal modifications required to adapt the PPO algorithm developed
 1489 previously for the SRBE to on for the SRABE.
 1490

$$\begin{aligned} 1491 \nabla_\theta \tilde{V}_{\text{RAA}}^{\pi_\theta}(s) &= \nabla_\theta \left(\sum_{a \in \mathcal{A}} \pi_\theta(a|s) \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s, a) \right) \\ 1492 &= \sum_{a \in \mathcal{A}} \left(\nabla_\theta \pi_\theta(a|s) \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s, a) \right. \\ 1493 &\quad \left. + \pi_\theta(a|s) \nabla_\theta \min \left\{ \max \left\{ \tilde{V}_{\text{RAA}}^\pi(f(s, a)), r_{\text{RAA}}(s) \right\}, q(s) \right\} \right) \\ 1494 &= \sum_{a \in \mathcal{A}} \left(\nabla_\theta \pi_\theta(a|s) \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s, a) \right. \\ 1495 &\quad \left. + \pi_\theta(a|s) \left[q(s) < \tilde{V}_{\text{RAA}}^\pi(f(s, a)) < r_{\text{RAA}}(s) \right] \nabla_\theta \tilde{V}_{\text{RAA}}^\pi(f(s, a)) \right) \end{aligned} \quad (37)$$

$$1501 = \sum_{s' \in \mathcal{S}} \left[\left(\sum_{k=0}^{\infty} \Pr(s \rightarrow s', k, \pi) \right) \sum_{a \in \mathcal{A}} \nabla_\theta \pi_\theta(a|s') \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s', a) \right] \quad (38)$$

$$\begin{aligned} 1502 &= \sum_{s' \in \mathcal{S}} \left[\left(\sum_{k=0}^{\infty} \Pr(s \rightarrow s', k, \pi) \right) \sum_{a \in \mathcal{A}} \pi_\theta(a|s') \frac{\nabla_\theta \pi_\theta(a|s')}{\pi_\theta(a|s')} \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s', a) \right] \\ 1503 &= \sum_{s' \in \mathcal{S}} \left[\left(\sum_{k=0}^{\infty} \Pr(s \rightarrow s', k, \pi) \right) \mathbb{E}_{a \sim \pi_\theta(s')} [\nabla_\theta \ln \pi_\theta(a|s') \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s', a)] \right] \\ 1504 &\propto \mathbb{E}_{s' \sim d'_\pi(s)} \mathbb{E}_{a \sim \pi_\theta(s')} [\nabla_\theta \ln \pi_\theta(a|s') \tilde{Q}_{\text{RAA}}^{\pi_\theta}(s', a)], \end{aligned}$$

1512 where the equality between (37) and (38) comes from rolling out the term $\nabla_\theta \tilde{V}_{\text{RAA}}^\pi(f(s, a))$ (see
 1513 Chapter 13.2 in Sutton and Barto (2018) for details), and where $\Pr(s \rightarrow s', k, \pi)$ is the probability
 1514 that under the policy π , the system is in state s' at time k given that it is in state s at time 0. \square
 1515

1516 Note, Proposition 1 is vital to updating the actor in Algorithm 1.
 1517

1518 E THE DO-HJ-PPO ALGORITHM

1520 In this section, we outline the details of our Actor-Critic algorithm DO-HJ-PPO beyond the details
 1521 given in Algorithm 1.
 1522

1523 Algorithm 1 : DO-HJ-PPO (Actor-Critic)

1525 **Require:** Composed and Decomposed Actor parameters θ and θ_i , Composed and Decomposed
 1526 Critic parameters ω and ω_i , GAE λ , learning rate β_k and discount factor γ . Let B^γ and B_i^γ
 1527 represent the Bellman update and decomposed Bellman update for the users choice of problem
 1528 (RR or RAA).
 1529 1: Define *Composed Actor* and Critic \tilde{Q}
 1530 2: Define *Decomposed Actor(s)* and Critic(s) \tilde{Q}_i
 1531 3: **for** $k = 0, 1, \dots$ **do**
 1532 4: **for** $t = 0$ to $T - 1$ **do**
 1533 5: Sample trajectories for $\tau_t : \{\hat{s}_t, a_t, \hat{s}_{t+1}\}$
 1534 6: Define $\tilde{\ell}(s_t)$ with Decomposed Critics $\tilde{Q}_i(s_t)$ (Theorems 1 & 2)
 1535 7: **Composed Critic update:**
 1536
$$\omega \leftarrow \omega - \beta_k \nabla_\omega \tilde{Q}(\tau_t) \cdot (\tilde{Q}(\tau_t) - B^\gamma[\tilde{Q}, \tilde{r}](\tau_t))$$

 1537
 1538 8: Compute Bellman-GAE A_{HJ}^λ with B^γ
 1539 9: (Standard) update Composed Actor
 1540 10: **Decomposed Critic update(s):**
 1541
$$\omega \leftarrow \omega - \beta_k \nabla_\omega \tilde{Q}_i(\tau_t) \cdot (\tilde{Q}_i(\tau_t) - B_i^\gamma[\tilde{Q}_i](\tau_t))$$

 1542
 1543 11: Compute Bellman-GAE A_i^λ with B_i^γ
 1544 12: (Standard) update Decomposed Actor(s)
 1545 13: **end for**
 1546 14: **end for**
 1547 15: **return** parameter θ, ω

1549 In Algorithm 1, the Bellman update $B^\gamma[\tilde{Q}, \tilde{r}]$ differs for the RAA task and RR task, and the $B_i^\gamma[\tilde{Q}]$
 1550 differs between the reach, avoid, and reach-avoid tasks.
 1551

1552 E.1 THE SPECIAL BELLMAN UPDATES AND THE CORRESPONDING GAES

1554 Akin to previous HJ-RL policy algorithms, namely RCPO Yu et al. (2022b), RESPO Ganai et al.
 1555 (2023) and RCPPO So et al. (2024), DO-HJ-PPO fundamentally depends on the discounted HJ
 1556 Bellman updates Fisac et al. (2019). To solve the RAA and RR problems with the special rewards
 1557 defined in Theorems 1 & 2, DO-HJ-PPO utilizes the Reach, Avoid and Reach-Avoid Bellman updates,
 1558 given by

$$B_R^\gamma[Q | r](s, a) = (1 - \gamma)r(s) + \gamma \max \{r(s), Q(s, a)\}, \quad (39)$$

$$B_A^\gamma[Q | q](s, a) = (1 - \gamma)q(s) + \gamma \min \{q(s), Q(s, a)\}, \quad (40)$$

$$B_{RA}^\gamma[Q | r, q](s, a) = (1 - \gamma) \min \{r(s), q(s)\} + \gamma \min \{q(s), \max \{r(s), Q(s, a)\}\}. \quad (41)$$

1563 To improve our algorithm, we incorporate the Generalized Advantage Estimate corresponding to
 1564 these Bellman equations in the updates of the Actors. As outlined in Section A of So et al. (2024), the
 1565 GAE may be defined with a reduction function corresponding to the appropriate Bellman function

which will be applied over a trajectory roll-out. We generalize the Reach GAE definition given in So et al. (2024) to propose a Reach-Avoid GAE (the Avoid GAE is simply the flip of the Reach GAE) as all will be used in DO-HJ-PPO algorithm for either RAA or RR problems. Consider a reduction function $\phi_{RA}^{(n)} : \mathbb{R}^n \rightarrow \mathbb{R}$, defined by

$$\phi_{RA}^{(n)}(x_1, x_2, x_3, \dots, x_{2n+1}) = \phi_{RA}^{(1)}(x_1, x_2, \phi_{RA}^{(n-1)}(x_3, \dots, x_{2n+1})), \quad (42)$$

$$\phi_{RA}^{(1)}(x, y, z) = (1 - \gamma) \min\{x, y\} + \gamma \min\{y, \max\{x, z\}\}. \quad (43)$$

The k -step Reach-Avoid Bellman advantage $A_{RA}^{\pi(k)}$ is then given by,

$$A_{RA}^{(k)}(s) = \phi_{RA}^{(n)} \left(r(s_t), q(s_t), \dots, r(s_{t+k-1}), q(s_{t+k-1}), V^{(s_{t+k})} \right) - V^{(s_{t+k})}. \quad (44)$$

We may then define the Reach-Avoid GAE A_{RA}^λ as the λ -weighted sum over the advantage functions

$$A_{RA}^\lambda(s) = \frac{1}{1 - \lambda} \sum_{k=1}^{\infty} \lambda^k A_{RA}^{(k)}(s) \quad (45)$$

which may be approximated over any finite trajectory sample. See So et al. (2024) for further details.

E.2 MODIFICATIONS FROM STANDARD PPO

To address the RAA and RR problems, DO-HJ-PPO introduces several key modifications to the standard PPO framework Schulman et al. (2017):

Additional actor and critic networks are introduced to represent the decomposed objectives.

Rather than learning the decomposed objectives separately from the composed objective, DO-HJ-PPO optimizes all objectives simultaneously. This design choice is motivated by two primary factors: (i) simplicity and minor computational speed-up, and (ii) coupling between the decomposed and composed objectives during learning.

The decomposed trajectories are initialized using states sampled from the composed trajectory, we refer to as *coupled resets*.

While it is possible to estimate the decomposed objectives independently—i.e., prior to solving the composed task—this approach might lead to inaccurate or irrelevant value estimates in on-policy settings. For example, in the RAA problem, the decomposed objective may prioritize avoiding penalties, while the composed task requires reaching a reward region without incurring penalties. In such a case, a decomposed policy trained in isolation might converge to an optimal strategy within a reward-irrelevant region, misaligned with the overall task. Empirically, we observe that omitting coupled resets causes DO-HJ-PPO to perform no better than standard baselines such as CPPO, whereas their inclusion significantly improves performance.

The special RAA and RR rewards are defined using the decomposed critic values and updated using their corresponding Bellman equations.

This procedure is directly derived from our theoretical results (Theorems 1 and 2), which establish the validity of using modified rewards within the respective RA and R Bellman frameworks. These rewards are used to compute the composed critic target as well as the actor’s GAE. In Algorithm 1, this process is reflected in the critic and actor updates corresponding to the composed objective.

F DDQN DEMONSTRATION

As described in the paper, we demonstrate the novel RAA and RR problems in a 2D Q -learning problem where the value function may be observed easily. We juxtapose these solitons with those of the previously studied RA and R problems which consider more simple objectives. To solve all values, we employ the standard Double-Deep Q learning approach (DDQN) Van Hasselt et al. (2016) with only the special Bellman updates.

1620 F.1 GRID-WORLD ENVIRONMENT
1621

1622 The environment is taken from Hsu et al. (2021) and consists of two dimensions, $s = (x, y)$, and
1623 three actions, $a \in \{\text{left}, \text{straight}, \text{right}\}$, which allow the agent to maneuver through the space. The
1624 deterministic dynamics of the environment are defined by constant upward flow such that,

$$1625 \quad f((x_i, y_i), a_i) = \begin{cases} (x_{i-1}, y_{i+1}) & a_i = \text{left} \\ (x_i, y_{i+1}) & a_i = \text{straight} \\ (x_{i+1}, y_{i+1}) & a_i = \text{right} \end{cases} \quad (46)$$

1629 and if the agent reaches the boundary of the space, defined by $x \geq |2|$, $y \leq -2$ and $y \geq 10$, the
1630 trajectory is terminated. The 2D space is divided into 80×120 cells which the agent traverses
1631 through.

1632 **In the RA and RAA experiments**, the reward function r is defined as the negative signed-distance
1633 function to a box with dimensions $(x_c, y_c, w, h) = (0, 4.5, 2, 1.5)$, and thus is negative iff the agent is
1634 outside of the box. The penalty function q is defined as the minimum of three (positive) signed distance
1635 functions for boxes defined at $(x_c, y_c, w, h) = (\pm 0.75, 3, 1, 1)$ and $(x_c, y_c, w, h) = (0, 6, 2.5, 1)$,
1636 and thus is positive iff the agent is outside of all boxes.

1637 **In the R and RR experiments**, one or two rewards are used. In the R experiment, the reward function
1638 r is defined as the maximum of two negative signed-distance function of boxes with dimensions
1639 $(x_c, y_c, w, h) = (\pm 1.25, 0, 0.5, 2)$, and thus is negative iff the agent is outside of both boxes. In the
1640 RR experiment, the rewards r_1 and r_2 are defined as the negative signed distance functions of the
1641 same two boxes independently, and thus are positive if the agent is in one box or the other respectively.

1643 F.2 DDQN DETAILS
1644

1645 As per our theoretical results in Theorems 1 and 2, we may now perform DDQN to solve the RAA
1646 and RR problems with solely the previously studied Bellman updates for the RA Hsu et al. (2021)
1647 and R problems Fisac et al. (2019). We compare these solutions with those corresponding to the
1648 RA and R problems *without* the special RAA and RR targets, and hence solve the previously posed
1649 problems. For all experiments, we employ the same adapted algorithm as in Hsu et al. (2021), with
1650 no modification of the hyper-parameters given in Table 1.

1651 Table 1: Hyperparameters for DDQN Grid World

1652 DDQN hyperparameters	1653 Values
1654 Network Architecture	MLP
1655 Numbers of Hidden Layers	2
1656 Units per Hidden Layer	100, 20
1657 Hidden Layer Activation Function	tanh
1658 Optimizer	Adam
1659 Discount factor γ	0.9999
1660 Learning rate	1e-3
1661 Replay Buffer Size	1e5 transitions
1662 Replay Batch Size	100
1663 Train-Collect Interval	10
1664 Max Updates	4e6

1665 G BASELINES
1666

1668 In both RAA and RR problems, we employ Constrained PPO (CPPO) Achiam et al. (2017a) as the
1669 major baseline as it can handle secondary objectives which are reformulated as constraints. The
1670 algorithm was not designed to minimize its constraints necessarily but may do so in attempting to
1671 satisfy them. As a novel direction in RL, few algorithms have been designed to optimize max/min
1672 accumulated costs and thus CPPO serves as the best proxy. Below we also include a naively
1673 decomposed STL algorithm to offer some insight into direct approaches to optimizing the max/min
accumulated reward.

1674
1675

G.1 CPPO BASELINES

1676
1677
1678
1679
1680
1681
1682

Although CPPO formulations do not directly consider dual-objective optimization, the secondary objective in RAA (avoid penalty) or overall objective in RR (reach both rewards) may be transformed into constraints to be satisfied of a surrogate problem. For the RAA problem, this may be defined as

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \max_{t' \leq t} r(s_t^{\pi}) \right] \quad \text{s.t.} \quad \min_t q(s_t^{\pi}) \geq 0. \quad (47)$$

1683

For the RR problem, one might propose that the fairest comparison would be to formulate the surrogate problem in the same fashion, with achievement of both costs as a constraint, such that

1684
1685
1686
1687
1688

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \min \left\{ \max_{t' \leq t} r_1(s_t^{\pi}), \max_{t' \leq t} r_2(s_t^{\pi}) \right\} \right] \quad \text{s.t.} \quad \min \left\{ \max_t r_1(s_t^{\pi}), \max_t r_2(s_t^{\pi}) \right\} \geq 0, \quad (48)$$

1689
1690
1691

which we define as variant 1 (CPPO-v1). Empirically, however, we found this formulation to be the poorest by far, perhaps due to the abundance of the non-smooth combinations. We thus also compare with more naive formulations which relax the outer minimizations to summation in the reward

1692
1693
1694
1695

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \max_{t' \leq t} r_1(s_t^{\pi}) + \max_{t' \leq t} r_2(s_t^{\pi}) \right] \quad \text{s.t.} \quad \min \left\{ \max_t r_1(s_t^{\pi}), \max_t r_2(s_t^{\pi}) \right\} \geq 0, \quad (49)$$

1696

which we define as variant 2 (CPPOv2), and additionally, in the constraint

1697
1698
1699

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \max_{t' \leq t} r_1(s_t^{\pi}) + \max_{t' \leq t} r_2(s_t^{\pi}) \right] \quad \text{s.t.} \quad \max_t r_1(s_t^{\pi}) + \max_t r_2(s_t^{\pi}) \geq 0, \quad (50)$$

1700
1701
1702

which we define as variant 3 (CPPOv3). This last approach, although naive and seemingly unfair, vastly outperforms the other variants in the RR problem.

1703
1704

G.2 STL BASELINES

1705
1706
1707
1708
1709
1710

In contrast with constrained optimization, one might also incorporate the STL methods, which in the current context simply decompose and optimize the independent objectives. For the RAA problem, the standard RA solution serves as a trivial STL baseline since we may attempt to continuously attempt to reach the solution while avoiding the obstacle. In the RR case, we define a decomposed STL baseline (DSTL) which naively solves both R problems, and selects the one with lower value to achieve first.

1711
1712
1713

H DETAILS OF RAA & RR EXPERIMENTS: HOPPER

1714
1715
1716
1717
1718

The Hopper environment is taken from Gym Brockman et al. (2016) and So et al. (2024). In both RAA and RR problems, we define rewards and penalties based on the position of the Hopper head, which we denote as (x, y) in this section.

In the RAA task, the reward is defined as

1719
1720

$$r(x, y) = \sqrt{\|x - 2\| + |y - 1.4|} - 0.1 \quad (51)$$

1721
1722
1723
1724
1725

to incentive the Hopper to reach its head to the position at $(x, y) = (2, 1.4)$. The penalty q is defined as the minimum of signed distance functions to a ceiling obstacle at $(1, 0)$, wall obstacles at $x > 2$ and $x < 0$ and a floor obstacle at $y < 0.5$. In order to safely arrive at high reward (and always avoid the obstacles), the Hopper thus must pass under the ceiling and not dive or fall over in the achievement of the target, as is the natural behavior.

1726
1727

In the RR task, the first reward is defined again as

$$r_1(x, y) = \sqrt{\|x - 2\| + |y - 1.4|} - 0.1 \quad (52)$$

1728 to incentive the Hopper to reach its head to the position at $(x, y) = (2, 1.4)$, and the second reward as
 1729

$$r_2(x, y) = \sqrt{\|x - 0\| + |y - 1.4|} - 0.1 \quad (53)$$

1730 to incentive the Hopper to reach its head to the position at $(x, y) = (0, 1.4)$. In order to achieve both
 1731 rewards, the Hopper must thus hop both forwards and backwards without crashing or diving.
 1732

1733 In all experiments, the Hopper is initialized in the default standing posture at a random $x \in [0, 2]$ so
 1734 as to learn a position-agnostic policy. The DO-HJ-PPO parameters used to train these problems can
 1735 be found in Table 2.
 1736

1737 Table 2: Hyperparameters for Hopper Learning

1738 Hyperparameters for DO-HJ-PPO	1739 Values
1740 Network Architecture	MLP
1741 Units per Hidden Layer	256
1742 Numbers of Hidden Layers	2
1743 Hidden Layer Activation Function	tanh
1744 Entropy coefficient	Linear Decay 1e-2 → 0
1745 Optimizer	Adam
1746 Discount factor γ	Linear Anneal 0.995 → 0.999
1747 GAE lambda parameter	0.95
1748 Clip Ratio	0.2
1749 Actor Learning rate	Linear Decay 3e-4 → 0
Reward/Cost Critic Learning rate	Linear Decay 3e-4 → 0
Add'l Hyperparameters for CPPO	
K_P	1
K_I	1e-4
K_D	1

1756 I DETAILS OF RAA & RR EXPERIMENTS: F16

1757 The F16 environment is taken from So et al. (2024), including a F16 fighter jet with a 26 dimensional
 1758 observation. The jet is limited to a flight corridor with up to 2000 relative position north (x_{PN}), 1200
 1759 relative altitude (x_H), and ± 500 relative position east (x_{PE}).
 1760

1761 In the RAA task, the reward is defined as
 1762

$$r(x, y) = \frac{1}{5}|x_{PN} - 1500| - 50 \quad (54)$$

1763 to incentivize the F16 to fly through the geofence defined by the vertical slice at 1500 relative position
 1764 north. The penalty q is defined as the minimum of signed distance functions to geofence (wall)
 1765 obstacles at $x_{PN} > 2000$ and $|x_{PE}| > 500$ and a floor obstacle at $x_H < 0$. In order to safely arrive
 1766 at high reward (and always avoid the obstacles), the F16 thus must fly through the target geofence
 1767 and then evade crashing into the wall directly in front of it.
 1768

1769 In the RR task, the rewards are defined as
 1770

$$r_1(x_{PN}, x_H) = \frac{1}{5}\sqrt{\|x_{PN} - 1250\| + |y - 850|} - 30 \quad (55)$$

1771 and
 1772

$$r_2(x_{PN}, x_H) = \frac{1}{5}\sqrt{\|x_{PN} - 1250\| + |y - 350|} - 30 \quad (56)$$

1773 to incentive the F16 to reach both low and high-altitude horizontal cylinders. In order to achieve both
 1774 rewards, the F16 must thus aggressively pitch, roll and yaw between the two targets.
 1775

1776 In all experiments, the F16 is initialized with position $x_{PN} \in [250, 750]$, $x_H \in [300, 900]$, $x_{PE} \in$
 1777 $[-250, 250]$ and velocity in $v \in [200, 450]$. Additionally, the roll, pitch, and yaw are initialized with
 1778 $\pm \pi/16$ to simulate a variety of approaches to the flight corridor. Further details can be found in
 1779 So et al. (2024). The DO-HJ-PPO parameters used to train these problems can be found in Table 3.
 1780

Table 3: Hyperparameters for F16 Learning

Hyperparameters for DO-HJ-PPO	Values
Network Architecture	MLP
Units per Hidden Layer	256
Numbers of Hidden Layers	2
Hidden Layer Activation Function	tanh
Entropy coefficient	Linear Decay $1e-2 \rightarrow 0$
Optimizer	Adam
Discount factor γ	Linear Anneal $0.995 \rightarrow 0.999$
GAE lambda parameter	0.95
Clip Ratio	0.2
Actor Learning rate	Linear Decay $1e-3 \rightarrow 0$
Reward/Cost Critic Learning rate	Linear Decay $1e-3 \rightarrow 0$
Add'l Hyperparameters for CPPO	
K_P	1
K_I	$1e-4$
K_D	1

J BROADER IMPACTS

This paper touches on advancing fundamental methods for Reinforcement Learning. In particular, this work falls into the class of methods designed for Safe Reinforcement Learning. Methods in this class are primarily intended to prevent undesirable behaviors in virtual or cyber-physical systems, such as preventing crashes involving self-driving vehicles or potentially even unacceptable speech among chatbots. It is an unfortunate truth that safe learning methods can be repurposed for unintended use cases, such as to prevent a malicious agent from being captured, but the authors do not foresee the balance of potential beneficial and malicious applications of this method to be any greater than other typical methods in Safe Reinforcement Learning.

K ACKNOWLEDGMENTS

This section has been redacted for the purpose of anonymous review.