

Supplementary Material

Contents

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

ACHIEVEMENT SPEED RESULTS FROM DO-HJ-PPO EXPERIMENTS

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

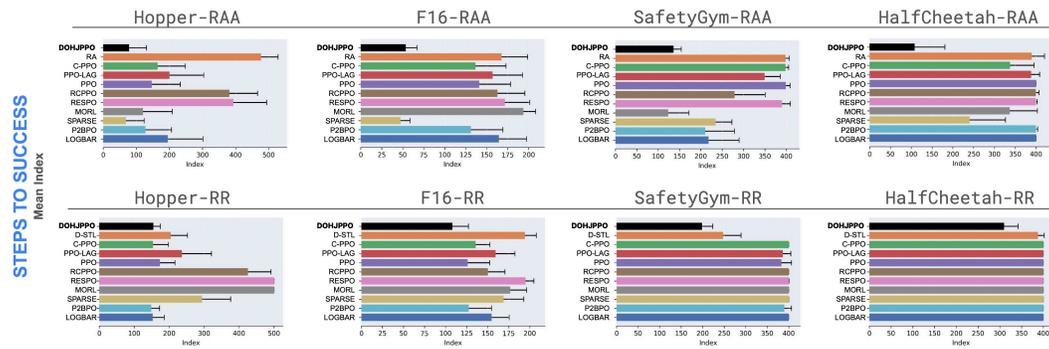


Figure 5: **Steps to Success** (\leftarrow) in RAA and RR Tasks for **DO-HJ-PPO** and Baselines For the same 1000 trajectories in Figure 4, we quantify here the number of steps until achievement of both tasks: reaching without crash afterward in the RAA, reaching both goal in the RR. **DO-HJ-PPO** is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

PROOF NOTATION

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

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 $[a, \cdot]$ 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

$$V_A^*(s) = \max_{\pi \in \Pi} \min_{\tau \in \mathbb{N}} q(\xi_s^\pi(\tau)),$$

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

$$V_{RAA}^*(s) = \max_{\pi \in \Pi} \min \left\{ \max_{\tau \in \mathbb{N}} r(\xi_s^\pi(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^\pi(\kappa)) \right\},$$

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:

$$v_A^*(s) = \max_{\mathbf{a} \in \mathbb{A}} \min_{\tau \in \mathbb{N}} q(\xi_s^\mathbf{a}(\tau)),$$

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

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

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*

$$811 \quad v_A^*(s) = \min_{\tau \in \mathbb{N}} q(\xi_s^\pi(\tau))$$

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

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

$$814 \quad \pi(s) \in \arg \max_{a \in \mathcal{A}} v_A^*(f(s, a)) \quad \forall s \in \mathcal{S}.$$

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

$$\begin{aligned} 816 \quad v_A^*(\xi_s^\pi(\tau + 1)) &= v_A^*(f(\xi_s^\pi(\tau), \pi(\xi_s^\pi(\tau)))) \\ 817 \quad &= \max_{a \in \mathcal{A}} v_A^*(f(\xi_s^\pi(\tau), a)) \\ 818 \quad &= \max_{a \in \mathcal{A}} \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q(\xi_{f(\xi_s^\pi(\tau), a)}^{\mathbf{a}}(\kappa)) \\ 819 \quad &= \max_{a \in \mathcal{A}} \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q(\xi_{\xi_s^\pi(\tau)}^{[\mathbf{a}, \mathbf{a}]}(\kappa + 1)) \\ 820 \quad &= \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q(\xi_{\xi_s^\pi(\tau)}^{\mathbf{a}}(\kappa + 1)) \\ 821 \quad &\geq \max_{\mathbf{a} \in \mathbb{A}} \min_{\kappa \in \mathbb{N}} q(\xi_{\xi_s^\pi(\tau)}^{\mathbf{a}}(\kappa)) \\ 822 \quad &\geq v_A^*(\xi_s^\pi(\tau)). \end{aligned}$$

823 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

$$824 \quad v_A^*(s) \geq \min_{\tau \in \mathbb{N}} q(\xi_s^\pi(\tau)) \geq \min_{\tau \in \mathbb{N}} v_A^*(\xi_s^\pi(\tau)) = v_A^*(\xi_s^\pi(0)) = v_A^*(s).$$

825 \square

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

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

$$828 \quad \tilde{v}_{\text{RA}}^*(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r_{\text{RAA}}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\}$$

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

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

$$831 \quad \max \emptyset = -\infty \quad \text{and} \quad \min \emptyset = +\infty.$$

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

$$833 \quad v^{\mathbf{a}}(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r_{\text{RAA}}(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}.$$

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

$$835 \quad V^\pi(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r_{\text{RAA}}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\}.$$

836 Then

$$837 \quad V^*(s) = \max_{\pi \in \Pi} \max_{\tau \in \mathbb{N}} \min \left\{ r_{\text{RAA}}(\xi_s^\pi(\tau)), \min_{\kappa \leq \tau} q(\xi_s^\pi(\kappa)) \right\} = \max_{\pi \in \Pi} V^\pi(s),$$

838 and

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

840 It is immediate that $v^*(s) \geq V^*(s)$ for each $s \in \mathcal{S}$, so it suffices to show the reverse inequality. Toward this end, it suffices to show that there is a $\pi \in \Pi$ for which $V^\pi(s) = v^*(s)$ for each $s \in \mathcal{S}$. 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 $s \in \bigcup_{n=0}^{\infty} S_n$. 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.

1. 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).

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

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

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

$$v^{\mathbf{a}}(s) = \min \left\{ \max \left\{ r_{\text{RAA}}(s), v^{\mathbf{a}1}(f(s, \mathbf{a}(0))) \right\}, q(s) \right\}, \quad \forall s \in \mathcal{S}, \mathbf{a} \in \mathbb{A},$$

from which it follows that

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

and

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

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

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

But then

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

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

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

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

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

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

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

From (17)-(20) we have

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

Now observe that

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

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

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

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

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

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

by (3), creating a contradiction.

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 $\xi_x^{\mathbf{a}}(\theta) \notin S_n$ 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). It follows from claim 1 that $\xi_x^{\mathbf{a}}(\theta) \in S_{n+1}$.

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

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

giving the desired contradiction.

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

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) and (10), we also have that $v_n^*(y) \geq \alpha_n \geq \alpha_{n+1}$. Together these establish (12) when $y \in S_n$.

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 possibilities. If there is some $a \in \mathcal{A}$ for which $f(y, a) \in S_n$, then (12) follows from claim 2. If 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 definition of v_0^*). Thus $\alpha_{n+1} = \min \{r_{\text{RAA}}(y), q(y)\}$ by (4), so

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

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

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

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

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

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

(\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 suprema over $\mathbf{a} \in \mathbb{A}$ and $\tau \in \mathbb{N}$ on both sides of this inequality gives the desired result.

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

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

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

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

$$1028 \quad \tilde{v}_{\text{RA}}^*(s) = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r_{\text{RAA}}(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} \quad (24)$$

$$1030 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), v_{\mathbb{A}}^*(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} \quad (25)$$

$$1032 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \min_{\kappa' \in \mathbb{N}} q(\xi_{\xi_s^{\mathbf{a}}(\tau)}^{\mathbf{a}'}(\kappa')), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$1034 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \min_{\kappa' \in \mathbb{N}} q(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \kappa')), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$1036 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa' \in \mathbb{N}} q(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \kappa')), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$1038 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \min \left\{ r(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau)), \min_{\kappa' \in \mathbb{N}} q(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\tau + \kappa')), \min_{\kappa \leq \tau} q(\xi_s^{[\mathbf{a}, \mathbf{a}']\tau}(\kappa)) \right\} \quad (26)$$

$$1040 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa' \in \mathbb{N}} q(\xi_s^{\mathbf{a}}(\tau + \kappa')), \min_{\kappa \leq \tau} q(\xi_s^{\mathbf{a}}(\kappa)) \right\} \quad (27)$$

$$1042 \quad = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$1044 \quad = \max_{\mathbf{a} \in \mathbb{A}} \min \left\{ \max_{\tau \in \mathbb{N}} r(\xi_s^{\mathbf{a}}(\tau)), \min_{\kappa \in \mathbb{N}} q(\xi_s^{\mathbf{a}}(\kappa)) \right\}$$

$$1046 \quad = v_{\text{RAA}}^*(s),$$

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

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

$$1056 \quad \mathcal{Y} = \{r(s) \mid s \in \mathcal{S}\} \quad \text{and} \quad \mathcal{Z} = \{q(s) \mid s \in \mathcal{S}\}.$$

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

$$1061 \quad \bar{\xi}_s^{\bar{\pi}}(t+1) = f(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\pi}(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\eta}_s^{\bar{\pi}}(t), \bar{\zeta}_s^{\bar{\pi}}(t))),$$

$$1062 \quad \bar{\eta}_s^{\bar{\pi}}(t+1) = \max \{r(\bar{\xi}_s^{\bar{\pi}}(t+1)), \bar{\eta}_s^{\bar{\pi}}(t)\},$$

$$1063 \quad \bar{\zeta}_s^{\bar{\pi}}(t+1) = \min \{q(\bar{\xi}_s^{\bar{\pi}}(t+1)), \bar{\zeta}_s^{\bar{\pi}}(t)\},$$

1064 for which $\bar{\xi}_s^{\bar{\pi}}(0) = s$, $\bar{\eta}_s^{\bar{\pi}}(0) = r(s)$, and $\bar{\zeta}_s^{\bar{\pi}}(0) = q(s)$.

1065 **Lemma 5.** *There is a $\bar{\pi} \in \bar{\Pi}$ such that*

$$1066 \quad v_{\text{RAA}}^*(s) = \min \left\{ \max_{\tau \in \mathbb{N}} r(\bar{\xi}_s^{\bar{\pi}}(\tau)), \min_{\tau \in \mathbb{N}} q(\bar{\xi}_s^{\bar{\pi}}(\tau)) \right\} \quad (28)$$

1067 for all $s \in \mathcal{S}$.

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

$$1074 \quad \tilde{v}_{\text{RA}}^*(s) = \max_{\tau \in \mathbb{N}} \min \left\{ r(\xi_s^{\pi}(\tau)), v_{\mathbb{A}}^*(\xi_s^{\pi}(\tau)), \min_{\kappa \leq \tau} q(\xi_s^{\pi}(\kappa)) \right\} \quad \forall s \in \mathcal{S},$$

$$1075 \quad v_{\mathbb{A}}^*(s) = \min_{\tau \in \mathbb{N}} q(\xi_s^{\theta}(\tau)) \quad \forall s \in \mathcal{S}.$$

1076 We introduce some useful notation we will use throughout the rest of the proof. For each $s \in \mathcal{S}$, let
1077 $[s]^+ = f(s, \pi(s))$, $[y]_s^+ = \max\{y, r([s]^+)\}$, $[z]_s^+ = \min\{z, q([s]^+)\}$.

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

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

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

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, $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \theta(\bar{x}_t)$, so that

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

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

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

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, we have that $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \pi(\bar{x}_t)$, so

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

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

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

There are again two cases:

($T < \infty$) 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 $t \leq T$. 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 by definition of T ,

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

and

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

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 $\min\{y_t^\circ, z_t^\circ, v_A^*(x_t^\circ)\}$ must achieve its maximal value at the smallest t for which it strictly decreases from t to $t+1$, i.e.

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

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

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

Thus

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

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

($T = \infty$) 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 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 definition of T ,

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

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
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

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

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

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

□

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

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

B PROOF OF RR MAIN THEOREM

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

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

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 rather than policies:

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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 $\tau_s = 0$. Then

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

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

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

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

$$\begin{aligned}
 1243 & \\
 1244 & v_{\mathbf{R}_1}^*(x) \geq v_{\mathbf{R}_1}^*(f(x, \pi_1(x))) \\
 1245 & \geq \max_{\mathbf{a} \in \mathbb{A}} r_1 \left(\xi_{f(x, \pi_1(x))}^{\mathbf{a}}(n) \right) \\
 1246 & \geq r_1 \left(\xi_{f(x, \pi_1(x))}^{\mathbf{a}_x|_1}(n) \right) \\
 1247 & \geq r_1 \left(\xi_x^{[\pi_1(x), \mathbf{a}_x|_1]}(n+1) \right) \\
 1248 & = r_1 \left(\xi_x^{[\pi_1(x), \mathbf{a}_x|_1]}(n+1) \right) \\
 1249 & = r_1 \left(\xi_x^{\mathbf{a}_x}(\tau_x) \right) \\
 1250 & = v_{\mathbf{R}_1}^*(x),
 \end{aligned}$$

1251 so that $v_{\mathbf{R}_1}^*(f(x, \pi_1(x))) = v_{\mathbf{R}_1}^*(x)$ and $\tau_{f(x, \pi_1(x))} \leq n$. It suffices to show

$$1252 \tau_{f(x, \pi_1(x))} = n, \quad (30)$$

1253 for then, by the induction assumption, we have

$$1254 r_1 \left(\xi_x^{\pi_1}(\tau_x) \right) = r_1 \left(\xi_{f(x, \pi_1(x))}^{\pi_1}(n) \right) = v_{\mathbf{R}_1}^*(f(x, \pi_1(x))) = v_{\mathbf{R}_1}^*(x).$$

1255 To show (30), assume instead that

$$1256 \tau_{f(x, \pi_1(x))} < n.$$

1257 But

$$\begin{aligned}
 1258 & v_{\mathbf{R}_1}^*(x) \geq \max_{\mathbf{a} \in \mathbb{A}} r_1 \left(\xi_x^{\mathbf{a}}(\tau_{f(x, \pi_1(x))} + 1) \right) \\
 1259 & \geq r_1 \left(\xi_x^{[\pi_1(x), \mathbf{a}_{f(x, \pi_1(x))}]}(\tau_{f(x, \pi_1(x))} + 1) \right) \\
 1260 & = r_1 \left(\xi_{f(x, \pi_1(x))}^{\mathbf{a}_{f(x, \pi_1(x))}}(\tau_{f(x, \pi_1(x))}) \right) \\
 1261 & = v_{\mathbf{R}_1}^*(f(x, \pi_1(x))) \\
 1262 & = v_{\mathbf{R}_1}^*(x),
 \end{aligned}$$

1263 so that

$$1264 v_{\mathbf{R}_1}^*(x) = \max_{\mathbf{a} \in \mathbb{A}} r_1 \left(\xi_x^{\mathbf{a}}(\tau_{f(x, \pi_1(x))} + 1) \right)$$

1265 and thus

$$1266 \tau_x \leq \tau_{f(x, \pi_1(x))} + 1 < n + 1,$$

1267 giving our desired contradiction. \square

1268 **Corollary 6.** For all $s \in \mathcal{S}$, we have $V_{\mathbf{R}_1}^*(s) = v_{\mathbf{R}_1}^*(s)$ and $V_{\mathbf{R}_2}^*(s) = v_{\mathbf{R}_2}^*(s)$.

1269 **Lemma 7.** There is a $\pi \in \Pi$ such that

$$1270 \tilde{v}_{\mathbf{R}}^*(s) = \max_{\tau \in \mathbb{N}} r_{\mathbf{R}\mathbf{R}} \left(\xi_s^{\pi}(\tau) \right).$$

1271 for all $s \in \mathcal{S}$.

1272 *Proof.* This lemma follows by precisely the same proof as the previous lemma, with r_1 , $v_{\mathbf{R}_1}^*$, and π_1 replaced with $r_{\mathbf{R}\mathbf{R}}$, $\tilde{v}_{\mathbf{R}}^*$, and π respectively. \square

1273 **Corollary 7.** For all $s \in \mathcal{S}$, we have $\tilde{V}_{\mathbf{R}}^*(s) = \tilde{v}_{\mathbf{R}}^*(s)$.

1274 **Lemma 8.** Let $\zeta_1 : \mathbb{N} \rightarrow \mathbb{R}$ and $\zeta_2 : \mathbb{N} \rightarrow \mathbb{R}$. Then

$$\begin{aligned}
 1275 & \sup_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}, \min \left\{ \sup_{\tau' \in \mathbb{N}} \zeta_1(\tau + \tau'), \zeta_2(\tau) \right\} \right\} \\
 1276 & = \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau), \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) \right\}.
 \end{aligned}$$

1277 *Proof.* We proceed by showing both inequalities corresponding to the above equality hold.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

(\leq) Observe that

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

(\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

$$\begin{aligned} & \sup_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\}, \min \left\{ \sup_{\tau' \in \mathbb{N}} \zeta_1(\tau + \tau'), \zeta_2(\tau) \right\} \right\} \\ & \geq \sup_{\tau \in \mathbb{N}} \min \left\{ \zeta_1(\tau), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau + \tau') \right\} \\ & \geq \min \left\{ \zeta_1(\tau_1), \sup_{\tau' \in \mathbb{N}} \zeta_2(\tau_1 + \tau') \right\} \\ & \geq \min \{ \zeta_1(\tau_1), \zeta_2(\tau_2) \} \\ & \geq \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau) - \varepsilon, \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) - \varepsilon \right\} \\ & = \min \left\{ \sup_{\tau \in \mathbb{N}} \zeta_1(\tau), \sup_{\tau \in \mathbb{N}} \zeta_2(\tau) \right\} - \varepsilon. \end{aligned}$$

But since $\varepsilon > 0$ was arbitrary, the desired inequality follows. □

Lemma 9. For each $s \in \mathcal{S}$,

$$\tilde{v}_R^*(s) = v_{RR}^*(s).$$

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

$$1351 \tilde{v}_{\text{R}}^*(s) = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} r_{\text{RR}}(\xi_s^{\mathbf{a}}(\tau)) \quad (31)$$

$$1352 = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), v_{\text{R}2}^*(\xi_s^{\mathbf{a}}(\tau)) \right\}, \min \left\{ v_{\text{R}1}^*(\xi_s^{\mathbf{a}}(\tau)), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \quad (32)$$

$$1353 = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_2(\xi_{\xi_s^{\mathbf{a}}(\tau)}^{\mathbf{a}'}(\tau')) \right\}, \right. \\ 1354 \left. \min \left\{ \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_1(\xi_{\xi_s^{\mathbf{a}}(\tau)}^{\mathbf{a}'}(\tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \\ 1355 = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau + \tau')) \right\}, \right. \\ 1356 \left. \min \left\{ \max_{\mathbf{a}' \in \mathbb{A}} \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \\ 1357 = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \max \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau + \tau')) \right\}, \right. \\ 1358 \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \\ 1359 = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max_{\mathbf{a}' \in \mathbb{A}} \max \left\{ \min \left\{ r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau + \tau')) \right\}, \right. \\ 1360 \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau + \tau')), r_2(\xi_s^{[\mathbf{a}, \mathbf{a}']_{\tau}}(\tau)) \right\} \right\} \quad (33)$$

$$1361 = \max_{\mathbf{a} \in \mathbb{A}} \max_{\tau \in \mathbb{N}} \max \left\{ \min \left\{ r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau' \in \mathbb{N}} r_2(\xi_s^{\mathbf{a}}(\tau + \tau')) \right\}, \right. \\ 1362 \left. \min \left\{ \max_{\tau' \in \mathbb{N}} r_1(\xi_s^{\mathbf{a}}(\tau + \tau')), r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \right\} \quad (34)$$

$$1363 = \max_{\mathbf{a} \in \mathbb{A}} \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\mathbf{a}}(\tau)), \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\mathbf{a}}(\tau)) \right\} \quad (35) \\ 1364 = v_{\text{RR}}^*(s),$$

1365 where the equality between [31](#) and [32](#) follows from Corollary [6](#), the equality between [33](#) and [34](#) follows from Lemma [3](#), and the equality between [34](#) and [35](#) follows from Lemma [8](#). \square

1366 Before the next lemma, we need to introduce two last pieces of notation. First, we let $\bar{\Pi}$ be the set of augmented policies $\bar{\pi} : \mathcal{S} \times \mathcal{Y} \times \mathcal{Z} \rightarrow \mathcal{A}$, as in the previous section, but where

$$1367 \mathcal{Y} = \{r_1(s) \mid s \in \mathcal{S}\} \quad \text{and} \quad \mathcal{Z} = \{r_2(s) \mid s \in \mathcal{S}\}.$$

1368 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 the solution of the evolution

$$1369 \bar{\xi}_s^{\bar{\pi}}(t+1) = f(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\pi}(\bar{\xi}_s^{\bar{\pi}}(t), \bar{\eta}_s^{\bar{\pi}}(t), \bar{\zeta}_s^{\bar{\pi}}(t))), \\ 1370 \bar{\eta}_s^{\bar{\pi}}(t+1) = \max \{r_1(\bar{\xi}_s^{\bar{\pi}}(t+1)), \bar{\eta}_s^{\bar{\pi}}(t)\}, \\ 1371 \bar{\zeta}_s^{\bar{\pi}}(t+1) = \max \{r_2(\bar{\xi}_s^{\bar{\pi}}(t+1)), \bar{\zeta}_s^{\bar{\pi}}(t)\},$$

1372 for which $\bar{\xi}_s^{\bar{\pi}}(0) = s$, $\bar{\eta}_s^{\bar{\pi}}(0) = r_1(s)$, and $\bar{\zeta}_s^{\bar{\pi}}(0) = r_2(s)$.

1373 **Lemma 10.** *There is a $\bar{\pi} \in \bar{\Pi}$ such that*

$$1374 v_{\text{RR}}^*(s) = \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\bar{\xi}_s^{\bar{\pi}}(\tau)), \max_{\tau \in \mathbb{N}} r_2(\bar{\xi}_s^{\bar{\pi}}(\tau)) \right\}$$

1375 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

$$1405 v_{R1}^*(s) = \max_{\tau \in \mathbb{N}} r_1(\xi_s^{\theta_1}(\tau)) \quad \forall s \in \mathcal{S},$$

$$1407 v_{R2}^*(s) = \max_{\tau \in \mathbb{N}} r_2(\xi_s^{\theta_2}(\tau)) \quad \forall s \in \mathcal{S},$$

$$1409 \tilde{v}_R^*(s) = \max_{\tau \in \mathbb{N}} \max \{ \min \{ r_1(\xi_s^\pi(\tau)), v_{R2}^*(\xi_s^\pi(\tau)) \}, \min \{ r_2(\xi_s^\pi(\tau)), v_{R1}^*(\xi_s^\pi(\tau)) \} \} \quad \forall s \in \mathcal{S}.$$

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

$$1412 \bar{\pi}(s, y, z) = \begin{cases} \pi(s) & \max\{y, z\} < \tilde{v}_R^*(s) \\ \theta_1(s) & \max\{y, z\} \geq \tilde{v}_R^*(s) \text{ and } y \leq z, \\ \theta_2(s) & \max\{y, z\} \geq \tilde{v}_R^*(s) \text{ and } y > z. \end{cases}$$

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

$$1419 v_{RR}^*(s) \leq \min \left\{ \max_{\tau \in \mathbb{N}} r_1(\bar{x}_\tau), \max_{\tau \in \mathbb{N}} r_2(\bar{x}_\tau) \right\}, \quad (36)$$

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

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

1425 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

$$\begin{aligned} 1426 \max_{t \in \mathbb{N}} \max \{ r_1(\bar{x}_t), r_2(\bar{x}_t) \} &< \max_{t \in \mathbb{N}} \tilde{v}_R^*(\bar{x}_t) \\ 1427 &= \tilde{v}_R^*(s) \\ 1428 &= \max_{\tau \in \mathbb{N}} \max \{ \min \{ r_1(x_\tau^\circ), v_{R2}^*(x_\tau^\circ) \}, \min \{ r_2(x_\tau^\circ), v_{R1}^*(x_\tau^\circ) \} \} \\ 1429 &= \max_{\tau \in \mathbb{N}} \max \{ \min \{ r_1(\bar{x}_\tau), v_{R2}^*(\bar{x}_\tau) \}, \min \{ r_2(\bar{x}_\tau), v_{R1}^*(\bar{x}_\tau) \} \} \\ 1430 &\leq \max_{\tau \in \mathbb{N}} \max \{ r_1(\bar{x}_\tau), r_2(\bar{x}_\tau) \}, \end{aligned}$$

1434 providing the desired contradiction.

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

$$1437 \max \{ r_1(\bar{x}_T), r_2(\bar{x}_T) \} \geq v_R^*(\bar{x}_T),$$

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

$$1440 \max \{ \min \{ r_1(x_{T'}^\circ), v_{R2}^*(x_{T'}^\circ) \}, \min \{ r_2(x_{T'}^\circ), v_{R1}^*(x_{T'}^\circ) \} \} = \tilde{v}_R^*(s),$$

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

$$1444 \max \{ \min \{ r_1(\bar{x}_{T'}), v_{R2}^*(\bar{x}_{T'}) \}, \min \{ r_2(\bar{x}_{T'}), v_{R1}^*(\bar{x}_{T'}) \} \} = \tilde{v}_R^*(s).$$

1446 But then

$$1447 \max \{ r_1(\bar{x}_{T'}), r_2(\bar{x}_{T'}) \} \geq \tilde{v}_R^*(s) \geq \tilde{v}_R^*(\bar{x}_{T'}).$$

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

- 1451 3. It follows from the previous step that

$$1452 \tilde{v}_R^*(\bar{x}_T) = \tilde{v}_R^*(x_{T'}^\circ) = \tilde{v}_R^*(s).$$

1454 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

$$1457 r_2(\bar{x}_t) < \tilde{v}_R^*(\bar{x}_T) \quad \forall t \in \mathbb{N}.$$

But then $\bar{\pi}(\bar{x}_t, \bar{y}_t, \bar{z}_t) = \theta_2(\bar{x}_t)$ for all $t \geq T$, so $v_{\mathbf{R}2}^*(\bar{x}_T) = \max_{t \geq T} r_2(\bar{x}_t) < \tilde{v}_{\mathbf{R}}^*(\bar{x}_T) \leq \tilde{v}_{\mathbf{R}}^*(s)$. Thus $r_2(x_{T'}^\circ) \leq v_{\mathbf{R}2}^*(x_{T'}^\circ) \leq v_{\mathbf{R}2}^*(x_T^\circ) = v_{\mathbf{R}2}^*(\bar{x}_T) < \tilde{v}_{\mathbf{R}}^*(s)$. It follows that

$$\max \{ \min \{ r_1(x_{T'}^\circ), v_{\mathbf{R}2}^*(x_{T'}^\circ) \}, \min \{ r_2(x_{T'}^\circ), v_{\mathbf{R}1}^*(x_{T'}^\circ) \} \} < \tilde{v}_{\mathbf{R}}^*(s),$$

contradicting our choice of T' .

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

□

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

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

C PROOF OF OPTIMALITY THEOREM

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

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

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

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

D THE SRABE AND ITS POLICY GRADIENT

Proof of Proposition 1 We here closely follow the proof of Theorem 3 in (4), which itself modifies the proofs of the Policy Gradient Theorems in Chapter 13.2 and 13.6 (60). We only make the minimal modifications required to adapt the PPO algorithm developed previously for the SRBE to on for the SRABE.

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

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

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

$$= \sum_{s' \in \mathcal{S}} \left[\left(\sum_{k=0}^{\infty} \Pr(s \rightarrow s', k, \pi) \right) \mathbb{E}_{a \sim \pi_\theta(s')} \left[\nabla_\theta \ln \pi_\theta(a|s') \tilde{Q}_{\mathbf{RAA}}^{\pi_\theta}(s', a) \right] \right]$$

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

where the equality between (37) and (38) comes from rolling out the term $\nabla_{\theta} \tilde{V}_{\text{RAA}}^{\pi}(f(s, a))$ (see Chapter 13.2 in (60) for details), and where $\text{Pr}(s \rightarrow s', k, \pi)$ is the probability 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

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

E THE DO-HJ-PPO ALGORITHM

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

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

Require: Composed and Decomposed Actor parameters θ and θ_i , Composed and Decomposed Critic parameters ω and ω_i , GAE λ , learning rate β_k and discount factor γ . Let B^{γ} and B_i^{γ} represent the Bellman update and decomposed Bellman update for the users choice of problem (RR or RAA).

- 1: Define *Composed* Actor and Critic \tilde{Q}
- 2: Define *Decomposed* Actor(s) and Critic(s) \tilde{Q}_i
- 3: **for** $k = 0, 1, \dots$ **do**
- 4: **for** $t = 0$ to $T - 1$ **do**
- 5: Sample trajectories for $\tau_t : \{\hat{s}_t, a_t, \hat{s}_{t+1}\}$
- 6: Define $\tilde{\ell}(s_t)$ with Decomposed Critics $\tilde{Q}_i(s_t)$ (Theorems 1 & 2)
- 7: **Composed Critic update:**

$$\omega \leftarrow \omega - \beta_k \nabla_{\omega} \tilde{Q}(\tau_t) \cdot (\tilde{Q}(\tau_t) - B^{\gamma}[\tilde{Q}, \tilde{r}](\tau_t))$$

- 8: Compute Bellman-GAE A_{HJ}^{λ} with B^{γ}
- 9: (Standard) update Composed Actor
- 10: **Decomposed Critic update(s):**

$$\omega \leftarrow \omega - \beta_k \nabla_{\omega} \tilde{Q}_i(\tau_t) \cdot (\tilde{Q}_i(\tau_t) - B_i^{\gamma}[\tilde{Q}_i](\tau_t))$$

- 11: Compute Bellman-GAE A_i^{λ} with B_i^{γ}
 - 12: (Standard) update Decomposed Actor(s)
 - 13: **end for**
 - 14: **end for**
 - 15: **return** parameter θ, ω
-

In Algorithm 1, the Bellman update $B^{\gamma}[\tilde{Q}, \tilde{r}]$ differs for the RAA task and RR task, and the $B_i^{\gamma}[\tilde{Q}_i]$ differs between the reach, avoid, and reach-avoid tasks. These Bellman updates are explicitly specified in the Supplementary Material.

E.1 THE SPECIAL BELLMAN UPDATES AND THE CORRESPONDING GAES

Akin to previous HJ-RL policy algorithms, namely RCPO (6), RESPO (3) and RCPPO (4), DO-HJ-PPO fundamentally depends on the discounted HJ Bellman updates (1). To solve the RAA and RR problems with the special rewards defined in Theorems 1 & 2, DO-HJ-PPO utilizes the Reach, Avoid and Reach-Avoid Bellman updates, given by

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

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

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

To improve our algorithm, we incorporate the Generalized Advantage Estimate corresponding to these Bellman equations in the updates of the Actors. As outlined in Section A of (4), the 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 (4) 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 (4) 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 (61):

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) (66) with only the special Bellman updates.

1620 F.1 GRID-WORLD ENVIRONMENT

1621
1622 The environment is taken from (2) and consists of two dimensions, $s = (x, y)$, and three actions,
1623 $a \in \{\text{left, straight, right}\}$, which allow the agent to maneuver through the space. The deterministic
1624 dynamics of the environment are defined by constant upward flow such that,

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

1628 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
1633 **In the RA and RAA experiments**, the reward function r is defined as the negative signed-distance
1634 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
1635 outside of the box. The penalty function q is defined as the minimum of three (positive) signed distance
1636 functions for boxes defined at $(x_c, y_c, w, h) = (\pm 0.75, 3, 1, 1)$ and $(x_c, y_c, w, h) = (0, 6, 2.5, 1)$,
1637 and thus is positive iff the agent is outside of all boxes.

1638 **In the R and RR experiments**, one or two rewards are used. In the R experiment, the reward function
1639 r is defined as the maximum of two negative signed-distance function of boxes with dimensions
1640 $(x_c, y_c, w, h) = (\pm 1.25, 0, 0.5, 2)$, and thus is negative iff the agent is outside of both boxes. In the
1641 RR experiment, the rewards r_1 and r_2 are defined as the negative signed distance functions of the
1642 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 (2) and R problems
1647 (1). We compare these solutions with those corresponding to the RA and R problems *without* the
1648 special RAA and RR targets, and hence solve the previously posed problems. For all experiments,
1649 we employ the same adapted algorithm as in (2), with no modification of the hyper-parameters given
1650 in Table 1.

1651
1652 Table 1: Hyperparameters for DDQN Grid World

DDQN hyperparameters	Values
Network Architecture	MLP
Numbers of Hidden Layers	2
Units per Hidden Layer	100, 20
Hidden Layer Activation Function	tanh
Optimizer	Adam
Discount factor γ	0.9999
Learning rate	1e-3
Replay Buffer Size	1e5 transitions
Replay Batch Size	100
Train-Collect Interval	10
Max Updates	4e6

1665 G BASELINES

1666
1667 In both RAA and RR problems, we employ Constrained PPO (CPPO) (8) as the major baseline as
1668 it can handle secondary objectives which are reformulated as constraints. The algorithm was not
1669 designed to minimize its constraints necessarily but may do so in attempting to satisfy them. As a
1670 novel direction in RL, few algorithms have been designed to optimize max/min accumulated costs
1671 and thus CPPO serves as the best proxy. Below we also include a naively decomposed STL algorithm
1672 to offer some insight into direct approaches to optimizing the max/min accumulated reward.
1673

1674 G.1 CPPO BASELINES

1675 Although CPPO formulations do not directly consider dual-objective optimization, the secondary
1676 objective in RAA (avoid penalty) or overall objective in RR (reach both rewards) may be transformed
1677 into constraints to be satisfied of a surrogate problem. For the RAA problem, this may be defined as
1678

$$1679 \max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \max_{t' \leq t} r(s_{t'}^{\pi}) \right] \quad \text{s.t.} \quad \min_t q(s_t^{\pi}) \geq 0. \quad (47)$$

1682 For the RR problem, one might propose that the fairest comparison would be to formulate the
1683 surrogate problem in the same fashion, with achievement of both costs as a constraint, such that
1684

$$1685 \max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \min \left\{ \max_{t' \leq t} r_1(s_{t'}^{\pi}), \max_{t' \leq t} r_2(s_{t'}^{\pi}) \right\} \right] \quad \text{s.t.} \quad \min \left\{ \max_t r_1(s_t^{\pi}), \max_t r_2(s_t^{\pi}) \right\} \geq 0, \quad (48)$$

1689 which we define as variant 1 (CPPO-v1). Empirically, however, we found this formulation to be the
1690 poorest by far, perhaps due to the abundance of the non-smooth combinations. We thus also compare
1691 with more naive formulations which relax the outer minimizations to summation in the reward
1692

$$1693 \max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \max_{t' \leq t} r_1(s_{t'}^{\pi}) + \max_{t' \leq t} r_2(s_{t'}^{\pi}) \right] \quad \text{s.t.} \quad \min \left\{ \max_t r_1(s_t^{\pi}), \max_t r_2(s_t^{\pi}) \right\} \geq 0, \quad (49)$$

1696 which we define as variant 2 (CPPOv2), and additionally, in the constraint
1697

$$1698 \max_{\pi} \mathbb{E}_{\pi} \left[\sum_t^{\infty} \gamma^t \max_{t' \leq t} r_1(s_{t'}^{\pi}) + \max_{t' \leq t} r_2(s_{t'}^{\pi}) \right] \quad \text{s.t.} \quad \max_t r_1(s_t^{\pi}) + \max_t r_2(s_t^{\pi}) \geq 0, \quad (50)$$

1700 which we define as variant 3 (CPPOv3). This last approach, although naive and seemingly unfair,
1701 vastly outperforms the other variants in the RR problem.
1702

1703 G.2 STL BASELINES

1705 In contrast with constrained optimization, one might also incorporate the STL methods, which in the
1706 current context simply decompose and optimize the independent objectives. For the RAA problem,
1707 the standard RA solution serves as a trivial STL baseline since we may attempt to continuously
1708 attempt to reach the solution while avoiding the obstacle. In the RR case, we define a decomposed
1709 STL baseline (DSTL) which naively solves both R problems, and selects the one with lower value to
1710 achieve first.
1711

1712 H DETAILS OF RAA & RR EXPERIMENTS: HOPPER

1713 The Hopper environment is taken from Gym (67) and (4). In both RAA and RR problems, we define
1714 rewards and penalties based on the position of the Hopper head, which we denote as (x, y) in this
1715 section.
1716

1717 In the RAA task, the reward is defined as
1718

$$1719 r(x, y) = \sqrt{\|x - 2\| + \|y - 1.4\|} - 0.1 \quad (51)$$

1720 to incentive the Hopper to reach its head to the position at $(x, y) = (2, 1.4)$. The penalty q is defined
1721 as the minimum of signed distance functions to a ceiling obstacle at $(1, 0)$, wall obstacles at $x > 2$
1722 and $x < 0$ and a floor obstacle at $y < 0.5$. In order to safely arrive at high reward (and always
1723 avoid the obstacles), the Hopper thus must pass under the ceiling and not dive or fall over in the
1724 achievement of the target, as is the natural behavior.
1725

1726 In the RR task, the first reward is defined again as
1727

$$1728 r_1(x, y) = \sqrt{\|x - 2\| + \|y - 1.4\|} - 0.1 \quad (52)$$

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

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

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

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

Table 2: Hyperparameters for Hopper 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 $3e-4 \rightarrow 0$
Reward/Cost Critic Learning rate	Linear Decay $3e-4 \rightarrow 0$
Add'l Hyperparameters for CPPO	
K_P	1
K_I	$1e-4$
K_D	1

I DETAILS OF RAA & RR EXPERIMENTS: F16

The F16 environment is taken from (4), including a F16 fighter jet with a 26 dimensional observation. The jet is limited to a flight corridor with up to 2000 relative position north (x_{PN}), 1200 relative altitude (x_H), and ± 500 relative position east (x_{PE}).

In the RAA task, the reward is defined as

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

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

In the RR task, the rewards are defined as

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

and

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

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

In all experiments, the F16 is initialized with position $x_{PN} \in [250, 750]$, $x_H \in [300, 900]$, $x_{PE} \in [-250, 250]$ and velocity in $v \in [200, 450]$. Additionally, the roll, pitch, and yaw are initialized with $\pm\pi/16$ to simulate a variety of approaches to the flight corridor. Further details can be found in (4). The DO-HJ-PPO parameters used to train these problems can be found in Table 3.

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.