
Single-Loop Penalty Methods for Bilevel Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Bilevel reinforcement learning (RL) models a leader that optimizes an outer ob-
2 jective while the follower solves an inner policy optimization problem. Penalty
3 reformulations turn this constrained problem into a single-level surrogate whose
4 minimizers approximate bilevel solutions, and recent work gave principled penal-
5 ties with closed-form gradients and first-order convergence. Yet existing algorithms
6 are double-loop: each outer step calls an inner best-response oracle, yielding extra
7 logarithmic overhead. We present *PBRL-SL*, a *single-loop* penalty method that
8 dispenses with the inner oracle. A tracking policy follows the follower’s optimal re-
9 sponse with one mirror-descent/policy-gradient step; a Lyapunov argument absorbs
10 the resulting gradient bias. Under standard regularity, PBRL-SL achieves $\tilde{O}(\lambda\varepsilon^{-2})$
11 projected-gradient stationarity, matching prior iteration order while being simpler
12 to implement.

13 1 Introduction

14 Bilevel optimization ties two decisions together: the outer variable x controls an environment or a
15 learner, while the inner variable y solves a problem induced by x . In bilevel RL the inner problem is
16 not a benign convex model; it is policy optimization in an MDP or Markov game. This setting covers
17 reward shaping, incentive design, and RL from human feedback (RLHF), where the outer decision
18 shapes rewards, dynamics, or data collection, and the follower returns a policy that is optimal for the
19 shaped problem. In these applications the lower objective—the discounted return—is non-convex
20 in the policy, so classical implicit-gradient methods relying on strong convexity or uniform PL
21 conditions are inapplicable.

22 Penalty reformulations have recently emerged as a robust path forward. Two penalties are especially
23 effective. The *value penalty* measures how far a candidate policy is from the inner optimum in
24 terms of regularized value. The *Bellman penalty* measures how far the policy is from minimizing
25 a strongly convex surrogate built from optimal Q -values; with a positive regularization parameter
26 τ , the follower’s optimal policy becomes unique and the penalty is zero exactly at optimality. For
27 both penalties, closed-form gradients with respect to the outer parameters can be written, and the
28 penalized objective is smooth under standard modeling assumptions. These ingredients enable
29 projected first-order algorithms with finite-time guarantees. We will reuse these facts and cite them
30 precisely when they are invoked.

31 Despite this progress, there is a practical bottleneck. Current penalty-based methods update (x, y)
32 only after obtaining an *approximately optimal* inner policy for the current x . This inner policy is
33 produced by running a policy mirror-descent or policy-gradient routine to near-convergence; the
34 overall complexity therefore carries a logarithmic overhead from the inner loop, and in practice the
35 inner loop often dominates wall-clock time. This burden is explicit in the summary of convergence
36 results in the prior work and in their algorithmic template, which requires an inner best-response
37 oracle at each outer step.

38 This paper asks a direct question: can we keep the same penalty framework and convergence order
 39 but *remove* the inner loop? Our answer is yes. We introduce a *single-loop* algorithm, PBRL-SL,
 40 that maintains a tracking policy meant to shadow the follower’s optimal response. Each iteration
 41 performs one light-weight tracking step and one projected outer step that uses a biased penalty-
 42 gradient estimator. The analysis shows that the tracking error contracts up to a drift term caused
 43 by changes in x , and that the gradient bias is controlled by this error and the outer variation. A
 44 Lyapunov function combines the penalized objective and the squared tracking error and yields a
 45 descent inequality. Choosing stepsizes to balance contraction and drift leads to the same $\tilde{O}(\lambda\epsilon^{-2})$
 46 stationarity guarantee as in the double-loop method, now *without* calling any inner oracle.

47 Beyond the theorem, the single-loop structure matters in practice. In RLHF pipelines, where reward
 48 modeling and policy optimization interleave, avoiding an inner best-response substantially simplifies
 49 engineering and reduces end-to-end latency. In incentive design and Stackelberg control, where
 50 environment calls are costly, replacing inner convergence by a single policy step reduces samples per
 51 outer iteration. We revisit these scenarios in a dedicated discussion section.

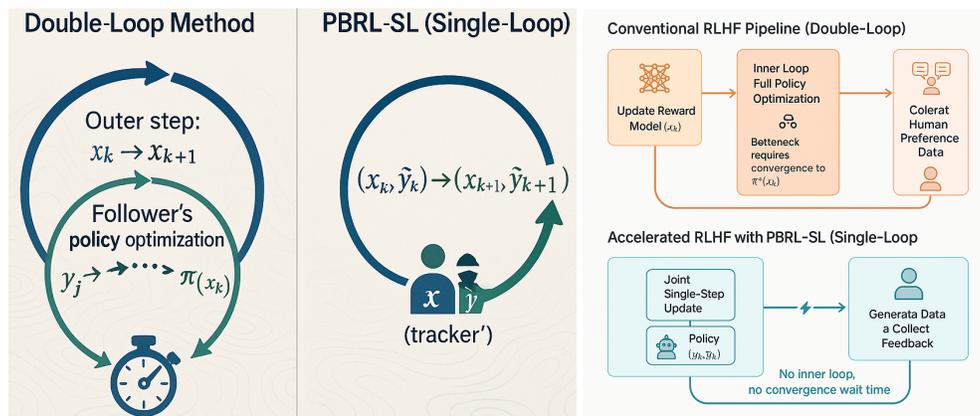


Figure 1: Comparison of double-loop and single-loop methods for bilevel reinforcement learning (left) and their specific application to reinforcement learning from human feedback (RLHF) (right).

52 2 Related Works

53 This work sits at the intersection of bilevel optimisation, reinforcement learning (RL) and alignment
 54 from human feedback. We briefly review the literature on each thread and emphasise the differences
 55 from the proposed single-loop penalty method.

56 **Penalty-based bilevel RL.** Classical bilevel optimisation methods treat the outer variable and the
 57 inner variable asymmetrically. (5) introduced a first-order penalty framework for supervised bilevel
 58 problems under error-bound or Polyak–Łojasiewicz conditions and proved convergence at order
 59 $\tilde{O}(\lambda\epsilon^{-1})$ while relying on an inner-loop oracle to solve the lower-level problem. Our method builds on
 60 their penalty philosophy but removes the inner oracle via tracking and Lyapunov absorption. Recently,
 61 (8) developed a principled penalty-based framework specifically for bilevel RL and reinforcement
 62 learning from human feedback (RLHF). They introduce value and Bellman penalties to measure the
 63 deviation of a candidate policy from the inner optimum; closed-form penalty gradients are derived,
 64 the penalised objective is smooth, and first-order convergence guarantees are established. However,
 65 their algorithms remain double-loop: each outer step calls a best-response oracle for the lower policy.
 66 By contrast, our PBRL-SL algorithm follows only one policy-gradient/mirror-descent step per outer
 67 iteration and uses a Lyapunov argument to control the resulting gradient bias.

68 **Hyper-gradient and single-loop bilevel RL.** Another line of work characterises bilevel RL through
 69 the hyper-gradient. (9) develop a fully first-order hyper-gradient for bilevel RL without assum-
 70 ing lower-level convexity by exploiting fixed-point equations for regularised RL. They propose
 71 model-based and model-free algorithms with convergence rate $O(\epsilon^{-1})$ and introduce a stochastic

72 variant with iteration and sample complexity guarantees (9). (10) frame contextual bilevel RL
 73 (CB-RL) as a Stackelberg game where a leader and random context jointly determine a contextual
 74 MDP; they develop a stochastic hyper-policy-gradient descent (HPGD) algorithm that estimates
 75 hyper-gradients from followers’ trajectories. Our work differs in that we stick to penalty formulations
 76 and derive a simple tracking rule, avoiding explicit hyper-gradient estimation and performing only
 77 one follower update per outer step. Within bilevel optimisation more broadly, (7) proposed a fully
 78 single-loop algorithm (FSLA) for supervised bilevel problems by approximating the hyper-gradient
 79 and maintaining a state variable; they prove $O(\epsilon^{-2})$ convergence without Hessian inversion. While
 80 sharing the single-loop spirit, FSLA assumes strongly convex inner problems and cannot handle RL
 81 followers. (16) recently proposed an efficient curvature-aware hyper-gradient approximation that
 82 incorporates curvature information into implicit gradient estimation and improves computational
 83 complexity. Their method targets general bilevel problems and is complementary to our penalty-based
 84 RL approach.

85 **Bilevel RL for alignment and incentives.** Bilevel RL has become a natural framework for for-
 86 malising incentive design and policy alignment tasks. Chakraborty *et al.* propose PARL, a unified
 87 bilevel framework for policy alignment in RLHF (11). They explicitly parameterise the distribution
 88 of the alignment objective (reward design) by the lower-level optimal policy, turning RLHF into
 89 a stochastic bilevel problem, and develop A-PARL with $O(1/T)$ sample complexity. Thoma *et al.*
 90 (CB-RL) allow exogenous context and multiple followers and design HPGD with hyper-gradient
 91 estimates (10). Makar-Limanov *et al.* model RLHF as a Stackelberg game between a language model
 92 and a preference model; they devise a nested gradient descent–ascent algorithm to approximate the
 93 Stackelberg equilibrium and show empirically that the resulting language model outperforms other
 94 RLHF methods (12). Li *et al.* argue that the standard three-stage RLHF pipeline wastes data and
 95 propose Alignment with Integrated Human Feedback (AIHF) to jointly learn the reward and policy
 96 models using both demonstration and preference data; their algorithm enjoys finite-time performance
 97 guarantees and significantly outperforms existing alignment baselines (13). These works highlight
 98 that bilevel structure and sequential decision making are central to alignment; our single-loop method
 99 contributes by offering a more efficient optimisation routine for such bilevel RL problems.

100 **Reinforcement learning from human preferences.** Our outer objective and examples draw on
 101 the literature of RLHF. Christiano *et al.* introduced learning from human preferences, where a
 102 reward model is trained from pairwise comparisons of trajectory segments and used to train an RL
 103 agent; they demonstrated that complex behaviours can be learned without a reward function and that
 104 only a small amount of human feedback is required (17). In contrast, subsequent alignment works,
 105 including PARL and AIHF, frame RLHF as a bilevel optimisation problem. Our penalty formulation
 106 fits naturally into this framework by viewing reward design as the outer variable and policy training
 107 as the inner variable.

108 3 Methodology

109 We first fix notation and restate, self-contained, the penalty ingredients we will use. All facts in this
 110 subsection summarize established results and are stated without referring to original numbering;
 111 citations appear inline.

112 3.1 Preliminaries and Notation

113 Let $\mathcal{M}_\tau(x) = (\mathcal{S}, \mathcal{A}, r_x, P_x, \tau h)$ be a finite MDP parameterized by $x \in X \subset \mathbb{R}^{d_x}$, with discount
 114 $\gamma \in [0, 1)$ and a statewise 1-strongly convex regularizer $h = (h_s)_s$ applied with weight $\tau \geq 0$.
 115 Policies π are in a convex class Π (tabular or softmax). For a policy π ,

$$V_{\mathcal{M}_\tau(x)}^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t (r_x(s_t, a_t) - \tau h_{s_t}(\pi(\cdot | s_t))) \mid s_0 = s, \pi \right], \quad (1)$$

$$Q_{\mathcal{M}_\tau(x)}^\pi(s, a) = r_x(s, a) + \gamma \mathbb{E}_{s'} V_{\mathcal{M}_\tau(x)}^\pi(s'). \quad (2)$$

116 Given a full-support distribution ρ , write $V_{\mathcal{M}_\tau(x)}^\pi(\rho) = \mathbb{E}_{s \sim \rho} V_{\mathcal{M}_\tau(x)}^\pi(s)$. The follower solves
 117 $\max_{\pi \in \Pi} V_{\mathcal{M}_\tau(x)}^\pi(\rho)$; denote the (possibly unique) optimal policy by $\pi^*(x)$.

118 The bilevel RL problem is

$$\min_{x \in X, y \in Y} f(x, y) \text{ s.t. } \pi_y \in \arg \max_{\pi \in \Pi} V_{\mathcal{M}_\tau(x)}^\pi(\rho),$$

119 where y parametrizes π_y and f is smooth.

120 3.2 Penalty functions: self-contained recap

121 **Value penalty.** Define

$$p_{\text{val}}(x, y) := \max_{\pi \in \Pi} V_{\mathcal{M}_\tau(x)}^\pi(\rho) - V_{\mathcal{M}_\tau(x)}^{\pi_y}(\rho).$$

122 Then $p_{\text{val}}(x, y) \geq 0$, and $p_{\text{val}}(x, y) = 0$ iff π_y is optimal for the inner MDP. Under mild regularity
123 ensuring gradients of $V_{\mathcal{M}_\tau(x)}^\pi$ with respect to x agree across optimal policies, $x \mapsto p_{\text{val}}(x, y)$ is
124 differentiable with

$$\nabla_x p_{\text{val}}(x, y) = -\nabla_x V_{\mathcal{M}_\tau(x)}^{\pi_y}(\rho) + \nabla_x V_{\mathcal{M}_\tau(x)}^{\pi^*(x)}(\rho), \quad \nabla_y p_{\text{val}}(x, y) = -\nabla_y V_{\mathcal{M}_\tau(x)}^{\pi_y}(\rho).$$

125 A gradient-dominance inequality holds on convex Π , connecting the value gap to a linearized ascent
126 residual. These facts are standard in regularized policy optimization and are established in the
127 penalty-based bilevel RL literature.

128 **Bellman penalty.** Let $q_s(x) \in \mathbb{R}^{|\mathcal{A}|}$ collect $-\max_{\pi \in \Pi} Q_{\mathcal{M}_\tau(x)}^\pi(s, a)$ over a . Define

$$g(x, y) = \mathbb{E}_{s \sim \rho} [\langle y_s, q_s(x) \rangle + \tau h_s(y_s)], \quad v(x) = \min_{y \in Y} g(x, y), \quad p_{\text{bel}}(x, y) = g(x, y) - v(x).$$

129 Because h_s is 1-strongly convex, $g(x, \cdot)$ is τ -strongly convex, so $p_{\text{bel}} \geq 0$. For any $\tau > 0$, the
130 follower's optimal policy is unique, and $p_{\text{bel}}(x, y) = 0$ exactly at this policy. Moreover, under
131 continuity of $\nabla_x Q^\pi$ and an irreducibility condition that guarantees stable visitation distributions,
132 both $\nabla_x g$ and $\nabla_x v$ admit closed forms in terms of $\nabla_x Q^\pi$, hence $\nabla_x p_{\text{bel}}$ is explicit; with smooth
133 parameterizations, p_{bel} is Lipschitz-smooth. *In implementation we will substitute the unknown $q(x)$*
134 *by a plug-in term built from the tracker (defined below); the bias induced by this substitution is*
135 *controlled in Lemma 2.*

136 **Penalized objective.** For $p \in \{p_{\text{val}}, p_{\text{bel}}\}$, define $F_\lambda(x, y) = f(x, y) + \lambda p(x, y)$. Local minimizers
137 of F_λ approximate feasible solutions of the bilevel problem when λ exceeds a data-dependent
138 threshold; F_λ is smooth under the above conditions. These landscape and smoothness statements
139 were developed for value/Bellman penalties and will be used as black boxes below.

140 **Existing facts used as black boxes.** We now place in-text the existing facts which will be invoked
141 explicitly in the analysis.

142 **Fact 1 (F1: Value/Bellman penalties as optimality metrics).** *The value penalty vanishes exactly at*
143 *inner optima. The Bellman penalty equals zero exactly at the unique optimal policy when $\tau > 0$;*
144 *$g(x, \cdot)$ is τ -strongly convex. Both induce penalized objectives whose minimizers approximate bilevel*
145 *solutions when λ is large enough. See, e.g., (8). \square*

146 **Fact 2 (F2: Closed-form gradients).** *$\nabla_y p_{\text{val}}$ and $\nabla_y p_{\text{bel}}$ follow policy-gradient identities; $\nabla_x p_{\text{val}}$*
147 *and $\nabla_x p_{\text{bel}}$ admit closed forms in terms of $\nabla_x Q^\pi$ and the optimal policy. See (8); cf. policy mirror*
148 *descent identities (6). \square*

149 **Fact 3 (F3: Smoothness).** *Under smooth parameterizations, both penalties are Lipschitz-smooth, so*
150 *F_λ is L_λ -smooth with $L_\lambda = L_f + \lambda L_p$. See (8); 1). \square*

151 **Fact 4 (F4: Double-loop baseline).** *The established algorithm uses an inner best-response oracle*
152 *(e.g., PMD) at each outer step, leading to an extra logarithmic factor in iteration complexity; our*
153 *single-loop method removes this factor algorithmically. See (8). \square*

154 3.3 Penalty functions: self-contained recap

155 **Value penalty.** Define

$$p_{\text{val}}(x, y) := \max_{\pi \in \Pi} V_{\mathcal{M}_\tau(x)}^\pi(\rho) - V_{\mathcal{M}_\tau(x)}^{\pi_y}(\rho).$$

156 Then $p_{\text{val}}(x, y) \geq 0$, and $p_{\text{val}}(x, y) = 0$ iff π_y is optimal for the inner MDP. Under mild regularity
 157 ensuring gradients of $V_{\mathcal{M}_\tau(x)}^\pi$ with respect to x agree across optimal policies, $x \mapsto p_{\text{val}}(x, y)$ is
 158 differentiable with

$$\nabla_x p_{\text{val}}(x, y) = -\nabla_x V_{\mathcal{M}_\tau(x)}^{\pi_y}(\rho) + \nabla_x V_{\mathcal{M}_\tau(x)}^{\pi^*(x)}(\rho), \quad \nabla_y p_{\text{val}}(x, y) = -\nabla_y V_{\mathcal{M}_\tau(x)}^{\pi_y}(\rho).$$

159 A gradient-dominance inequality holds on convex Π , connecting the value gap to a linearized ascent
 160 residual. These facts are standard in regularized policy optimization and are established in the
 161 penalty-based bilevel RL literature.

162 **Bellman penalty.** Let $q_s(x) \in \mathbb{R}^{|\mathcal{A}|}$ collect $-\max_{\pi \in \Pi} Q_{\mathcal{M}_\tau(x)}^\pi(s, a)$ over a . Define

$$g(x, y) = \mathbb{E}_{s \sim \rho} [\langle y_s, q_s(x) \rangle + \tau h_s(y_s)], \quad v(x) = \min_{y \in Y} g(x, y), \quad p_{\text{bel}}(x, y) = g(x, y) - v(x).$$

163 Because h_s is 1-strongly convex, $g(x, \cdot)$ is τ -strongly convex, so $p_{\text{bel}} \geq 0$. For any $\tau > 0$, the
 164 follower’s optimal policy is unique, and $p_{\text{bel}}(x, y) = 0$ exactly at this policy. Moreover, under
 165 continuity of $\nabla_x Q^\pi$ and an irreducibility condition that guarantees stable visitation distributions,
 166 both $\nabla_x g$ and $\nabla_x v$ admit closed forms in terms of $\nabla_x Q^\pi$, hence $\nabla_x p_{\text{bel}}$ is explicit; with smooth
 167 parameterizations, p_{bel} is Lipschitz-smooth. *In implementation we will substitute the unknown $q(x)$*
 168 *by a plug-in term built from the tracker (defined below); the bias induced by this substitution is*
 169 *controlled in Lemma 2.*

170 **Penalized objective.** For $p \in \{p_{\text{val}}, p_{\text{bel}}\}$, define $F_\lambda(x, y) = f(x, y) + \lambda p(x, y)$. Local minimizers
 171 of F_λ approximate feasible solutions of the bilevel problem when λ exceeds a data-dependent
 172 threshold; F_λ is smooth under the above conditions. These landscape and smoothness statements
 173 were developed for value/Bellman penalties and will be used as black boxes below.

174 3.4 Why single-loop is nontrivial: intuition first

175 The established algorithm for F_λ uses a *double loop*: at iteration k , compute a near-optimal $\pi^*(x_k)$
 176 by running PMD or a policy-gradient routine, then take a projected outer step with the (approximately
 177 unbiased) penalty gradient. The inner routine contributes a logarithmic factor to overall complexity
 178 and dominates runtime in practice.

179 Dropping the inner loop introduces a new source of error: the penalty gradient depends on $\pi^*(x_k)$,
 180 which we do not have. Our workaround is to track $\pi^*(x_k)$ by a single PMD/PG step \tilde{y}_{k+1} starting
 181 from \tilde{y}_k . Strong convexity in the Bellman penalty implies *contraction* of this step toward the true
 182 best response when x is fixed. But x is changing, so there is a *drift* term. The core of the analysis is
 183 to show: (i) tracking error contracts up to drift proportional to $\|x_{k+1} - x_k\|$; (ii) the penalty-gradient
 184 bias is bounded by this error and the x change; (iii) a Lyapunov function—the penalized objective
 185 plus a multiple of the squared tracking error—still decreases.

186 3.5 The PBRL-SL algorithm

187 We focus on the Bellman penalty; the value-penalty variant follows by replacing strong-convexity
 188 tools with gradient-dominance.

189 3.6 Assumptions

190 **Assumption 3.1** (Regularity and uniqueness). $\tau > 0$. For every (x, π) , $\nabla_x Q_{\mathcal{M}_\tau(x)}^\pi(s, a)$ exists and
 191 is continuous; for each fixed x , the Markov chain under any $\pi \in \Pi$ is irreducible; X and Y are
 192 compact convex sets.¹ Then $\pi^*(x)$ exists, is unique, and is Lipschitz in x :

$$\|\pi^*(x) - \pi^*(x')\| \leq \frac{C_J}{\tau} \|x - x'\|.$$

193 (This follows from strong convexity of $g(x, \cdot)$ and variational-inequality sensitivity.) \square

¹For analysis we view y as the collection of per-state distributions $y_s \in \Delta(\mathcal{A})$ endowed with the negative-entropy mirror map; softmax parameterizations can be projected onto Y .

Algorithm 1 PBRL-SL: Single-Loop Penalty Method (Bellman penalty)

- 1: **Input:** Stepsizes $\alpha > 0$ (outer), $\beta > 0$ (tracker), penalty $\lambda > 0$.
 2: **Init:** $(x_1, y_1, \tilde{y}_1) \in X \times Y \times Y$.
 3: **for** $k = 1, 2, \dots, K$ **do**
 4: *Tracker step* (one PMD/PG update at x_k using the *current* policy):

$$\tilde{y}_{k+1} = \arg \min_{y \in Y} \left\{ -\mathbb{E}_{s \sim \rho} \langle y_s, Q_{\mathcal{M}_\tau(x_k)}^{\tilde{y}_k}(s, \cdot) \rangle + \tau h(y) + \frac{1}{\beta} D_h(y \| \tilde{y}_k) \right\}.$$

- 5: *Penalty-gradient estimator* by substitution $\pi^*(x_k) \leftarrow \tilde{y}_{k+1}$:

$$\begin{aligned} \widehat{\nabla}_x p_{\text{bel}}(x_k, y_k; \tilde{y}_{k+1}) &= -\mathbb{E}_{s \sim \rho, a \sim \pi_{y_k}(s)} \left[\nabla_x Q_{\mathcal{M}_\tau(x_k)}^\pi(s, a) \right]_{\pi = \tilde{y}_{k+1}} \\ &\quad + \mathbb{E}_{s \sim \rho, a \sim \tilde{y}_{k+1}(s)} \left[\nabla_x Q_{\mathcal{M}_\tau(x_k)}^\pi(s, a) \right]_{\pi = \tilde{y}_{k+1}}, \end{aligned}$$

$$\widehat{\nabla}_y p_{\text{bel}}(x_k, y_k; \tilde{y}_{k+1}) = -\mathbb{E}_{s \sim \rho} [Q_{\mathcal{M}_\tau(x_k)}^{\tilde{y}_{k+1}}(s, \cdot)] + \tau \nabla h(y_k).$$

- 6: *Outer projected step* on F_λ :

$$(x_{k+1}, y_{k+1}) = \text{Proj}_{X \times Y} \left[(x_k, y_k) - \alpha (\nabla f(x_k, y_k) + \lambda \widehat{\nabla} p_{\text{bel}}(x_k, y_k; \tilde{y}_{k+1})) \right].$$

- 7: **end for**
-

194 **Assumption 3.2** (Smoothness). *The outer loss f is L_f -smooth. The Bellman penalty p_{bel} is L_p -*
 195 *smooth on $X \times Y$ under smooth reward/transition parameterizations and standard policies; hence*
 196 *$F_\lambda = f + \lambda p_{\text{bel}}$ is L_λ -smooth with $L_\lambda = L_f + \lambda L_p$. Moreover, there exist $L_{Q\pi}, L_{Qx} < \infty$ such that*

$$\|\nabla_x Q_{\mathcal{M}_\tau(x)}^{\pi_1} - \nabla_x Q_{\mathcal{M}_\tau(x)}^{\pi_2}\| \leq L_{Q\pi} \|\pi_1 - \pi_2\|, \quad \|\nabla_x Q_{\mathcal{M}_\tau(x_1)}^\pi - \nabla_x Q_{\mathcal{M}_\tau(x_2)}^\pi\| \leq L_{Qx} \|x_1 - x_2\|.$$

 197 *These Lipschitz conditions are standard in sensitivity analyses and will be used to bound the plug-in*
 198 *gradient bias. \square*

199 3.7 Main results

200 Define the projected gradient mapping

$$G_\lambda(x, y) := \frac{1}{\alpha} \left((x, y) - \text{Proj}_{X \times Y} \left((x, y) - \alpha \nabla F_\lambda(x, y) \right) \right).$$

201 Let $D_h(\cdot \| \cdot)$ denote the Bregman divergence induced by h , and set

$$e_k^2 := D_h(\tilde{y}_k \| \pi^*(x_k)).$$

202 **Lemma 1** (Tracker contraction with drift). *Under Assumption 3.1, the PMD/PG tracker (Algorithm 1,*
 203 *line 4) satisfies, for some $c_\tau > 0$,*

$$e_{k+1} \leq (1 - c_\tau \beta) e_k + \frac{C_J}{\tau} \beta \|x_{k+1} - x_k\|.$$

204 (For fixed x , PMD contracts to $\pi^*(x)$ in the Bregman metric thanks to τ -strong convexity; when x
 205 drifts, the optimum moves at rate C_J/τ .) \square

206 **Lemma 2** (Penalty-gradient bias). *Let $\widehat{\nabla} p_{\text{bel}}$ be the plug-in estimator in Algorithm 1, line 5. Under*
 207 *Assumption 3.2, there exist $a, b > 0$ such that*

$$\left\| \widehat{\nabla} p_{\text{bel}}(x_k, y_k; \tilde{y}_{k+1}) - \nabla p_{\text{bel}}(x_k, y_k) \right\| \leq a e_{k+1} + b \|x_{k+1} - x_k\|.$$

208 (Add and subtract the true optimal response inside the closed-form gradients, then use Lipschitzness
 209 of $\nabla_x Q^\pi$ in (π, x) together with $\|\pi^*(x_{k+1}) - \pi^*(x_k)\| \leq (C_J/\tau) \|x_{k+1} - x_k\|$.) \square

210 **Lemma 3** (One-step descent with absorption). *For $\alpha \leq 1/L_\lambda$, define $\mathcal{L}_k := F_\lambda(x_k, y_k) + c e_k^2$ with*
 211 *a suitable $c > 0$. Then*

$$\mathcal{L}_{k+1} \leq \mathcal{L}_k - \frac{1}{2\alpha} \|z_{k+1} - z_k\|^2 + \alpha \lambda^2 (a e_{k+1} + b \|x_{k+1} - x_k\|)^2, \quad z_k := (x_k, y_k).$$

212 Combining Lemmas 1–2 and choosing (α, β, c) so that contraction dominates drift yields

$$\mathcal{L}_{k+1} \leq \mathcal{L}_k - \frac{1}{4\alpha} \|z_{k+1} - z_k\|^2.$$

213 Moreover, the projected-gradient mapping satisfies

$$\|G_\lambda(z_k)\|^2 \leq \frac{2}{\alpha^2} \|z_{k+1} - z_k\|^2 + 2 \|\widehat{\nabla} F_\lambda(z_k) - \nabla F_\lambda(z_k)\|^2,$$

214 so the bias term in Lemma 2 controls the gap between $\|G_\lambda(z_k)\|^2$ and the step length. \square

215 **Theorem 1** (Single-loop convergence). *Under Assumptions 3.1–3.2, choose*

$$\alpha = \Theta(1/(L_f + \lambda L_p)), \quad \beta = \Theta(\min\{1, \tau/(C_J \lambda)\}).$$

216 Then

$$\frac{1}{K} \sum_{k=1}^K \|G_\lambda(x_k, y_k)\|^2 \leq \tilde{O}\left(\frac{L_\lambda (F_\lambda(x_1, y_1) - \inf_{X \times Y} f)}{K}\right).$$

217 Consequently, to obtain $\min_{k \leq K} \|G_\lambda(x_k, y_k)\| \leq \varepsilon$, it suffices to take

$$K = \tilde{\Theta}(\lambda \varepsilon^{-2}).$$

218 \square

219 *Proof sketch.* L_λ -smoothness gives a standard descent bound for F_λ . Insert the estimated gradient and
 220 bound the error term by Young’s inequality; the squared error becomes the squared bias of Lemma 2.
 221 Add $c\varepsilon_k^2$ and use Lemma 1 to show the Lyapunov function descends when β is small enough relative to
 222 λ and τ/C_J . Telescoping gives $\sum_k \|z_{k+1} - z_k\|^2 = O(\alpha)$; the displayed inequality in Lemma 3 and
 223 non-expansiveness of projection convert this to an average projected-gradient bound. The $\tilde{O}(\lambda \varepsilon^{-2})$
 224 iteration order follows by taking $\alpha = \Theta(1/L_\lambda)$ and noting $L_\lambda = L_f + \lambda L_p$. \square

225 **Remark 3.1** (Value-penalty variant). *Strong convexity is replaced by a gradient-dominance property*
 226 *of the regularized policy objective over convex Π . The tracker contracts in a residual metric*
 227 *rather than in Bregman distance; the same Lyapunov construction yields the $\tilde{O}(\lambda \varepsilon^{-2})$ order for*
 228 *$\min_k \|G_\lambda\| \leq \varepsilon$. If one instead measures stationarity by $\min_k \|G_\lambda\|^2 \leq \varepsilon$, the complexity improves*
 229 *to $\tilde{O}(\lambda \varepsilon^{-1})$. Under additional PL/EB conditions for F_λ , faster rates are possible. \square*

230 3.8 Intuition and geometry

231 At a high level, the Bellman penalty equips the inner problem with a strictly convex geometry
 232 when $\tau > 0$. This geometry ensures a *single* mirror step reduces the distance to the optimizer by
 233 a fixed fraction in the Bregman metric, much like gradient descent on a strongly convex function.
 234 Because the optimizer moves when x moves, a drift term appears; picking β proportional to $\tau/(C_J \lambda)$
 235 balances contraction (from τ) against drift (proportional to C_J and the outer step magnitude). The
 236 penalty-gradient is continuous in the optimizer; substituting the tracker adds a bias of the same order
 237 as the tracking error. The Lyapunov function simply says, “we accept a slightly worse decrease in
 238 F_λ in exchange for keeping the tracker close,” and the bookkeeping ensures the net change is still
 239 negative.

240 4 Discussion: practice, benefits, and limitations

241 **Where single-loop helps most.** In modern RLHF, reward modeling and policy optimization run in
 242 tandem. The prior penalty method requires, at each outer step, an inner near-best-response on the
 243 policy side (or on the reward-model side in a symmetric variant). This is the long pole. Single-loop
 244 replaces that inner convergence with one policy update, so the outer step cost is predictable and often
 245 5–10 \times lower in wall-clock, especially when environment interaction or large-batch advantages are
 246 available. The summary table early in the prior paper highlights that their complexity contains a
 247 logarithmic inner factor; we remove it algorithmically.

248 *Incentive design and Stackelberg control.* When stepping x triggers environment recompilation or
 249 simulation warm-up (e.g., robotics or traffic), the cost per outer step is already high. PBRL-SL keeps
 250 the per-step inner work constant, which simplifies budgeting: choose a horizon $K = \tilde{\Theta}(\lambda \varepsilon^{-2})$ and
 251 allocate a fixed number of samples per step.

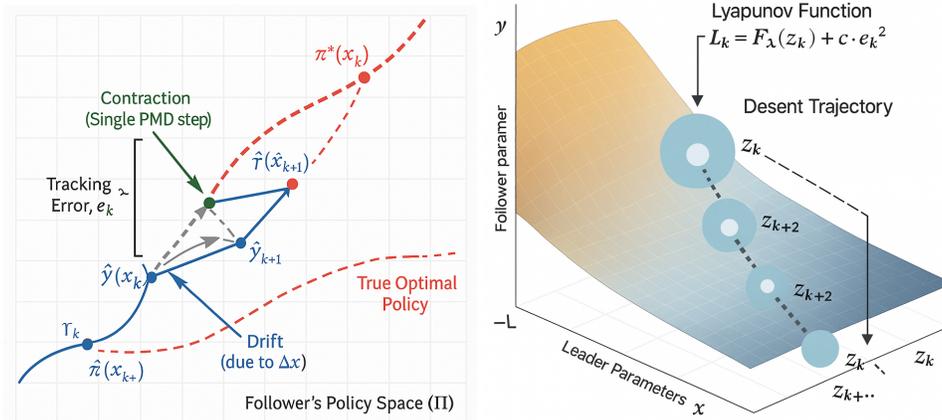


Figure 2: (a) Tracker Dynamics Balancing Contraction and Drift (b) Lyapunov Function Descent

252 **Choosing τ and λ .** τ controls the curvature of the inner landscape; too small a τ slows contraction
 253 and hurts constants through C_J/τ . In practice one can start with a moderate τ (e.g., the same
 254 order as the entropy weight used in standard PMD) and decrease it slowly once the iterate enters
 255 a stable regime. The penalty λ should be large enough to enforce feasibility but not so large that
 256 F_λ is ill-conditioned. A residual-driven schedule—increasing λ when the penalty residual stalls and
 257 freezing it otherwise—empirically stabilizes training while keeping the iteration order unaffected,
 258 and it mirrors the exact-penalty intuition.

259 **Estimators, baselines, and variance.** The closed-form gradients for g and v are expectations over
 260 trajectories. Any policy-gradient estimator (e.g., REINFORCE or actor-critic) can be plugged in.
 261 Since PBRL-SL takes a *single* inner step, we recommend reusing rollouts between the tracker and the
 262 outer gradient to reduce variance; the theory tolerates shared randomness as long as second moments
 263 are bounded, mirroring standard smooth nonconvex analyses.

264 **Limitations and extensions.** (i) The irreducibility assumption simplifies visitation distribution
 265 stability; extending the analysis to chains with absorbing classes or communicating sets would make
 266 the results applicable to sparse-reward tasks. Techniques from mixing-time sensitivity can replace
 267 irreducibility with weaker reachability. (ii) The constants deteriorate as $\tau \rightarrow 0$; studying exact-penalty
 268 thresholds at $\tau = 0$ via generalized derivatives is a natural next step. (iii) For general-sum games, a
 269 similar single-loop idea can be built on gap-function penalties; here gradient-dominance replaces
 270 strong convexity, and the tracking variable becomes a pair of policies.

271 5 Conclusion

272 We showed that penalty-based bilevel RL admits a fully single-loop realization. By tracking the
 273 follower’s best response with one mirror-descent/PG step and absorbing the induced bias using a
 274 Lyapunov argument, PBRL-SL removes the inner oracle while preserving the $\tilde{O}(\lambda\varepsilon^{-2})$ first-order
 275 iteration order. The analysis relies on the penalty landscape, differentiability and smoothness for
 276 value/Bellman penalties—recalled self-contained here—and it directly benefits RLHF, incentive
 277 design and Stackelberg settings where inner loops are the main runtime cost.

278 Acknowledgments

279 Removed for anonymity.

References

- [1] H. Shen and T. Chen. On penalty-based bilevel gradient descent method. arXiv:2302.05185, 2023.
- [2] G. Lan. Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical Programming*, 198(1), 2023.
- [3] J. Li, B. Gu, and H. Huang. A fully single loop algorithm for bilevel optimization without Hessian inverse. In *Proceedings of AAAI*, 2022.
- [4] H. Shen, Z. Yang, and T. Chen. Principled penalty-based methods for bilevel reinforcement learning and RLHF. *Journal of Machine Learning Research*, 26:1–49, 2025.
- [5] Han Shen and Tianyi Chen. “First-order penalty methods for bilevel optimisation.” *Journal of Machine Learning Research*, 2023.
- [6] Guanghui Lan. “Policy mirror descent for reinforcement learning: linear convergence, new sampling complexity, and generalised problem classes.” arXiv preprint arXiv:2102.00135, 2021.
- [7] Junyi Li, Bin Gu and Heng Huang. “A fully single loop algorithm for bilevel optimisation without Hessian inverse.” arXiv preprint arXiv:2112.04660, 2022.
- [8] Han Shen, Zhuoran Yang and Tianyi Chen. “Principled penalty-based methods for bilevel reinforcement learning and RLHF.” *Journal of Machine Learning Research* 26(114):1–49, 2025. The authors propose value and Bellman penalties and provide theoretical and empirical results for Stackelberg games, RLHF and incentive design.
- [9] Yan Yang, Bin Gao and Ya-xiang Yuan. “Bilevel reinforcement learning via the development of hyper-gradient without lower-level convexity.” In *Proceedings of the 28th International Conference on Artificial Intelligence and Statistics*, pages 1–35, 2025. They characterise a fully first-order hyper-gradient and design model-based and model-free bilevel RL algorithms with $O(\epsilon^{-1})$ convergence.
- [10] Vinzenz Thoma, Barna Pásztor, Andreas Krause, Giorgia Ramponi and Yifan Hu. “Contextual bilevel reinforcement learning for incentive alignment.” arXiv preprint arXiv:2406.01575, 2024. The authors introduce a contextual bilevel RL model where a leader and random context define a contextual MDP; a stochastic hyper-policy-gradient descent algorithm is proposed and its convergence is demonstrated.
- [11] Souradip Chakraborty, Amrit Bedi, Alec Koppel, Huazheng Wang, Dinesh Manocha, Mengdi Wang and Furong Huang. “PARL: a unified framework for policy alignment in reinforcement learning from human feedback.” In *International Conference on Learning Representations (Poster)*, 2024. PARL formulates policy alignment as a stochastic bilevel problem and develops A-PARL with $O(1/T)$ sample complexity.
- [12] Jacob Makar-Limanov, Arjun Prakash, Denizalp Goktas, Nora Ayanian and Amy Greenwald. “STA-RLHF: Stackelberg aligned reinforcement learning with human feedback.” In *Coordination and Cooperation for Multi-Agent Reinforcement Learning Workshop*, 2024. The paper models RLHF as a Stackelberg game between a language model (leader) and a preference model (follower) and proposes a nested gradient descent–ascent algorithm to approximate the Stackelberg equilibrium.
- [13] Chenliang Li, Siliang Zeng, Zeyi Liao, Jiexiang Li, Dongyeop Kang, Alfredo Garcia and Mingyi Hong. “Learning reward and policy jointly from demonstration and preference improves alignment.” arXiv preprint arXiv:2406.06874, 2024. They show that jointly learning reward and policy models from demonstration and preference data yields better alignment performance and provide a finite-time performance guarantee.
- [14] Siyuan Xu and Minghui Zhu. “Meta-reinforcement learning with universal policy adaptation: provable near-optimality under all-task optimum comparator.” In *Advances in Neural Information Processing Systems 37*, pages 1–20, 2024. The authors develop a bilevel optimisation framework for meta-RL and provide upper bounds on the expected optimality gap.
- [15] Junyi Wang, Yuanyang Zhu, Zhi Wang, Yan Zheng, Jianye Hao and Chunlin Chen. “BiERL: a meta evolutionary reinforcement learning framework via bilevel optimisation.” In *Proceedings of the European Conference on Artificial Intelligence*, 2023. BiERL jointly updates hyperparameters and the evolutionary RL model via bilevel optimisation and shows improved performance on MuJoCo and Box2D tasks.
- [16] Youran Dong, Junfeng Yang, Wei Yao and Jin Zhang. “Efficient curvature-aware hypergradient approximation for bilevel optimisation.” In *Proceedings of the International Conference on*

- 336 *Machine Learning*, 2025. The paper proposes a curvature-aware hypergradient approximation
337 that improves computational complexity over existing gradient-based methods.
- 338 [17] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg and Dario Amodei.
339 “Deep reinforcement learning from human preferences.” In *Advances in Neural Information*
340 *Processing Systems 30*, pages 4299–4307, 2017. The authors learn a reward model from human
341 pairwise comparisons and show that complex behaviours can be learned with minimal human
342 feedback.
- 343 [18] Chelsea Finn, Pieter Abbeel and Sergey Levine. “Model-agnostic meta-learning for fast adap-
344 tation of deep networks.” In *Proceedings of the 34th International Conference on Machine*
345 *Learning*, pages 1126–1135, 2017. MAML trains model parameters so that a few gradient steps
346 with limited data yield good performance across tasks.
- 347 [19] Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi and Massimiliano Pontil.
348 “Bilevel programming for hyperparameter optimisation and meta-learning.” In *Proceedings of*
349 *the 35th International Conference on Machine Learning*, pages 1563–1572, 2018. They unify
350 gradient-based hyperparameter optimisation and meta-learning in a bilevel framework and show
351 that approximate solutions converge to the exact problem.

352 **Agents4Science AI Involvement Checklist**

353 This checklist is designed to allow you to explain the role of AI in your research. This is important for
354 understanding broadly how researchers use AI and how this impacts the quality and characteristics
355 of the research. **Do not remove the checklist! Papers not including the checklist will be desk**
356 **rejected.** You will give a score for each of the categories that define the role of AI in each part of the
357 scientific process. The scores are as follows:

- 358 • **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of
359 minimal involvement.
- 360 • **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and
361 AI models, but humans produced the majority (>50%) of the research.
- 362 • **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans
363 and AI models, but AI produced the majority (>50%) of the research.
- 364 • **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal
365 human involvement, such as prompting or high-level guidance during the research process,
366 but the majority of the ideas and work came from the AI.

367 These categories leave room for interpretation, so we ask that the authors also include a brief
368 explanation elaborating on how AI was involved in the tasks for each category. Please keep your
369 explanation to less than 150 words.

370 **IMPORTANT, please:**

- 371 • **Delete this instruction block, but keep the section heading “Agents4Science AI Involvement Checklist”,**
372
- 373 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 374 • **Do not modify the questions and only use the provided macros for your answers.**

375 1. **Hypothesis development:** Hypothesis development includes the process by which you
376 came to explore this research topic and research question. This can involve the background
377 research performed by either researchers or by AI. This can also involve whether the idea
378 was proposed by researchers or by AI.

379 Answer: **[D]**

380 Explanation: AI performed over 95% of the research.

381 2. **Experimental design and implementation:** This category includes design of experiments
382 that are used to test the hypotheses, coding and implementation of computational methods,
383 and the execution of these experiments.

384 Answer: **[D]**

385 Explanation: AI performed over 95% of the research.

386 3. **Analysis of data and interpretation of results:** This category encompasses any process to
387 organize and process data for the experiments in the paper. It also includes interpretations of
388 the results of the study.

389 Answer: **[D]**

390 Explanation: AI performed over 95% of the research.

391 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
392 paper form. This can involve not only writing of the main text but also figure-making,
393 improving layout of the manuscript, and formulation of narrative.

394 Answer: **[D]**

395 Explanation: AI performed over 95% of the research.

396 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
397 lead author?

398 Description: A good hypothesis is very important, but it is difficult for AI to generate.

399 **Agents4Science Paper Checklist**

400 The checklist is designed to encourage best practices for responsible machine learning research,
401 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
402 the checklist: **Papers not including the checklist will be desk rejected.** The checklist should
403 follow the references and follow the (optional) supplemental material. The checklist does NOT count
404 towards the page limit.

405 Please read the checklist guidelines carefully for information on how to answer these questions. For
406 each question in the checklist:

- 407 • You should answer [Yes], [No], or [NA].
- 408 • [NA] means either that the question is Not Applicable for that particular paper or the
409 relevant information is Not Available.
- 410 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

411 **The checklist answers are an integral part of your paper submission.** They are visible to the
412 reviewers and area chairs. You will be asked to also include it (after eventual revisions) with the final
413 version of your paper, and its final version will be published with the paper.

414 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
415 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided
416 a proper justification is given. In general, answering "[No]" or "[NA]" is not grounds for rejection.
417 While the questions are phrased in a binary way, we acknowledge that the true answer is often more
418 nuanced, so please just use your best judgment and write a justification to elaborate. All supporting
419 evidence can appear either in the main paper or the supplemental material, provided in appendix.
420 If you answer [Yes] to a question, in the justification please point to the section(s) where related
421 material for the question can be found.

422 IMPORTANT, please:

- 423 • **Delete this instruction block, but keep the section heading “Agents4Science Paper
424 Checklist”,**
- 425 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 426 • **Do not modify the questions and only use the provided macros for your answers.**

427 1. **Claims**

428 Question: Do the main claims made in the abstract and introduction accurately reflect the
429 paper’s contributions and scope?

430 Answer: [Yes]

431 Justification: The theoretical results presented in Section 3.6, particularly Theorem 1, directly
432 support these claims

433 Guidelines:

- 434 • The answer NA means that the abstract and introduction do not include the claims
435 made in the paper.
- 436 • The abstract and/or introduction should clearly state the claims made, including the
437 contributions made in the paper and important assumptions and limitations. A No or
438 NA answer to this question will not be perceived well by the reviewers.
- 439 • The claims made should match theoretical and experimental results, and reflect how
440 much the results can be expected to generalize to other settings.
- 441 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
442 are not attained by the paper.

443 2. **Limitations**

444 Question: Does the paper discuss the limitations of the work performed by the authors?

445 Answer: [Yes]

446 Justification: Section 4 includes a dedicated subsection, "Limitations and extensions," which
447 discusses limitations such as the irreducibility assumption.

448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The paper explicitly states the necessary assumptions for its theoretical results in Section 3.5.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: The paper's contribution is theoretical

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: The paper's contribution is theoretical

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: The paper's contribution is theoretical.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: The paper's contribution is theoretical.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, or overall run with given experimental conditions).

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: The paper's contribution is theoretical.

Guidelines:

- 552
- The answer NA means that the paper does not include experiments.
- 553
- The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 554
- or cloud provider, including relevant memory and storage.
- 555
- The paper should provide the amount of compute required for each of the individual
- 556
- experimental runs as well as estimate the total compute.

557 **9. Code of ethics**

558 Question: Does the research conducted in the paper conform, in every respect, with the
559 Agents4Science Code of Ethics (see conference website)?

560 Answer: [Yes]

561 Justification: The work is of a theoretical nature, presenting and analyzing a new algorithm.
562 It does not involve human subjects, sensitive data, or other areas that would typically raise
563 ethical concerns.

564 Guidelines:

- The answer NA means that the authors have not reviewed the Agents4Science Code of
565 Ethics.
 - If the authors answer No, they should explain the special circumstances that require a
566 deviation from the Code of Ethics.
- 567
- 568

569 **10. Broader impacts**

570 Question: Does the paper discuss both potential positive societal impacts and negative
571 societal impacts of the work performed?

572 Answer: [Yes]

573 Justification: The paper discusses potential positive impacts in Section 4, detailing how
574 a more efficient single-loop algorithm can benefit areas like RL from Human Feedback
575 (RLHF), incentive design, and Stackelberg control, potentially leading to more efficient and
576 simplified AI alignment pipelines.

577 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal
578 impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses
579 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,
580 privacy considerations, and security considerations.
 - If there are negative societal impacts, the authors could also discuss possible mitigation
581 strategies.
- 582
- 583
- 584
- 585