

Algorithms with Smoothed Polynomial-Time Complexity for Deterministic Discounted-sum and Mean-payoff Games

Ali Asadi
IST Austria
ali.asadi@ist.ac.at

Krishnendu Chatterjee
IST Austria
kchatterjee@ist.ac.at

Ruichen Luo
IST Austria
rluo@ist.ac.at

Abstract

Turn-based deterministic games on graphs are two-player zero-sum games played on finite directed graphs, where the vertex set is partitioned between the Max and the adversarial Min Player. The respective player chooses outgoing edges from their vertices. We consider these games with the classical discounted-sum and mean-payoff objectives. The goal of the Max Player is to maximize the objective against the adversarial Min Player. These games lie in $NP \cap coNP$ and are among the rare combinatorial problems that belong to this complexity class, yet the existence of a polynomial-time algorithm is a major open question. All known deterministic algorithms require exponential time in the worst case, while the classical strategy iteration (policy iteration) algorithm is remarkably efficient in practice, which motivates the study of the smoothed complexity of these games. Previous results established a smoothed polynomial-time bound only for the restricted class of ergodic games and presented only a high-probability guarantee. Our main result is a new analysis technique for a variant of the strategy iteration algorithm, and we show these algorithms achieve smoothed polynomial complexity for general game graphs with both high-probability and expected runtime guarantees. In particular, our result removes the ergodicity restriction of prior work and strengthens the high-probability guarantee to the stronger expected-runtime guarantee, thereby resolving the open question on the smoothed complexity of these games.¹

¹Through personal communication, we are aware of independent and concurrent work by Bruno Loff and Mateusz Skomra, titled “Smoothed analysis of policy iteration and value iteration for deterministic discounted and mean-payoff games played on any graph”.

1 Introduction

Turn-based game graphs. The class of two-player turn-based (perfect-information) deterministic games is a widely studied theoretical model. The model consists of a finite directed graph $G = (V, E)$ where the vertices are partitioned into Player-Max and Player-Min vertices, and the corresponding player chooses outgoing edges of the vertices in control. This model extends classical graphs and is used to model interaction between two adversarial agents, and has many applications, e.g., in analysis of reactive synthesis (interaction of controller and environment) [PR89, RW87]; in adversarial planning [MB85, HZ98]. Moreover, this model closely corresponds to the notion of alternation in Turing machines which has been studied in [CKS81].

Discounted-sum and mean-payoff objectives. In the study of games over graphs, the two fundamental payoff functions or objectives are discounted-sum and mean-payoff. Every edge is assigned a real-valued reward, which represents the per-stage reward. A play (or infinite-walk) of the graph is an infinite-sequence of vertices, with edges between the consecutive vertices. The payoff of a play under discounted-sum (resp., mean-payoff) is the discounted-sum (resp., the long-run average) of the rewards that appear in the play. These objectives are the most well-studied and basic objectives studied in games over graphs [Sha53, Gil57, FV96, Put94], e.g., mean-payoff and discounted-sum games have been studied for quantitative reactive synthesis [DAH03, CPV15, BCH09].

Algorithmic results. Besides the practical motivation to study turn-based deterministic games with discounted-sum and mean-payoff objectives, they represent an intriguing class of problems for algorithmic study. The main computational problem is related to the optimal value for an objective, where the Max Player aims to maximize and the Min Player aims to minimize the payoff function. The decision problems for these games lie in $\text{NP} \cap \text{coNP}$ (also $\text{UP} \cap \text{coUP}$) [Con92, ZP96], but the existence of a polynomial-time algorithm is a long-standing and major open problem. All deterministic algorithms for this problem have a worst-case exponential-time complexity.

Smoothed complexity. While the worst-case complexity for the computational problem for these games is exponential, a classical algorithm (namely the policy-iteration algorithm) works very well in practice. Hence it is natural to consider the smoothed complexity of this problem [ST04], which is a classical complexity analysis between worst-case and average-case analysis. The smoothed complexity analysis of these problems has an interesting history: First, a polynomial-time smoothed complexity result for mean-payoff games was claimed in [BEF⁺11]; unfortunately this claim was incorrect and later retracted. Second, this result was partially recovered in [LS24] which shows smoothed polynomial-time complexity for the special class of ergodic games (where all vertices have the same value) with a high-probability guarantee.

Open problem. The class of ergodic games is often simpler to analyze, e.g., partially-observable MDPs (POMDPs) with ergodicity are decidable [CLS25] whereas in general they are undecidable [MHC03]; and for mean-payoff objectives the algorithmic study for ergodic games is much simpler [CLJ14, HK66]. Hence the first key open problem is to extend the smoothed polynomial time result from ergodic to general game graphs. Moreover, the previous result only presents a high-probability guarantee, whereas for smoothed complexity analysis a stronger and more desired guarantee is the expected runtime guarantee (see Remark 1). Thus, the second interesting open problem is to establish polynomial-time smoothed complexity with expected runtime guarantee for turn-based deterministic discounted-sum and mean-payoff games.

Our contribution. We answer the above open questions in the affirmative. Moreover, we show that variants of the policy-iteration algorithm achieve the polynomial-time smoothed complexity on general game graphs with expected runtime guarantee. In other words, we resolve the open question with variants of a well-known algorithm that also works well in practice.

68 *Related works.* The algorithmic study of discounted-sum and mean-payoff games has a long his-
69 tory [EM79, GKK88, ZP96, ACSSU24, DKZ19, Koz21]. All these deterministic algorithms have
70 worst-case exponential-time complexity for the general case, and even the known randomized al-
71 gorithms achieve only sub-exponential time [Lud95, BV07]. Moreover, discounted-sum games are
72 closely related to turn-based stochastic games, which subsume MDPs as a special case. Thus, it is
73 natural to ask whether our smoothed polynomial-time result extends to this broader class. Quite
74 interestingly, [CY23] establishes a super-polynomial-time lower bound for the smoothed complexity
75 of policy iteration on MDPs, which rules out such an extension to turn-based stochastic games.

76 2 Preliminaries

77 For ease of reference, Table 1 in Appendix A collects the recurring notation of the paper.

78 **Definition 1** (Game graphs). A game graph $G = (V_{\min}, V_{\max}, E)$ is a directed finite graph con-
79 sisting of two disjoint finite sets of vertices V_{\min} and V_{\max} , with vertex set $V \triangleq V_{\min} \uplus V_{\max}$ and
80 edges $E \subseteq V \times V$. We denote by $n \triangleq |V|$ and $m \triangleq |E|$. An edge $a = (v, v') \in E$ indicates the
81 transition from vertex v to v' . Without loss of generality, we consider that for every vertex $v \in V$,
82 there exists at least one $v' \in V$ such that $(v, v') \in E$.

83 **Dynamics and paths.** At the beginning of the game, a token is placed at an initial vertex
84 $v_0 \in V$. At any time t ($t \in \mathbb{Z}_{\geq 0}$): if $v_t \in V_{\min}$, then the Min Player chooses an edge $a_t \in E$, where
85 $a_t = (v_t, v_{t+1})$, and the token is moved from v_t to v_{t+1} ; and similarly, if $v_t \in V_{\max}$, then the Max
86 Player chooses an edge and moves the token. A *path* (or *play*) starting from an initial vertex v_0 is
87 an infinite sequence of vertices $\rho = v_0, v_1, v_2, \dots$ such that $(v_t, v_{t+1}) \in E$ for all $t \in \mathbb{Z}_{\geq 0}$.

88 **Objectives.** We consider *discounted-sum* and *mean-payoff* objectives, and in both cases there
89 is a reward vector $\mathbf{r} \in \mathbb{R}^E$ that assigns to each edge $a = (v, v') \in E$ a real number $\mathbf{r}(v, v')$ as its
90 per-stage reward.

- 91 • *Discounted-sum objective:* For the discounted-sum objective, a real discount factor $\gamma \in [0, 1)$
92 is given, and the game is denoted by $G^{\gamma, \mathbf{r}}$. For any path $\rho = v_0, v_1, v_2, \dots$, the discounted-sum
93 payoff is the infinite discounted sum of the rewards of the edges that appear in the path:

$$u^{\gamma, \mathbf{r}}(\rho) = \sum_{t=0}^{\infty} \gamma^t \mathbf{r}(v_t, v_{t+1}).$$

94 Since $\gamma \in [0, 1)$ and the graph is finite (i.e., rewards are bounded), this infinite sum converges.

- 95 • *Mean-payoff objectives.* For mean-payoff objectives, the same graph and reward vector are
96 considered, denoted by $G^{\text{mp}, \mathbf{r}}$, but the payoff function evaluates a path by its long-run average
97 reward instead of discounted sum. For any path $\rho = v_0, v_1, \dots$, the *mean-payoff* is:

$$u^{\text{mp}, \mathbf{r}}(\rho) = \liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{r}(v_t, v_{t+1}).$$

98 **Strategies.** A *strategy* (or *policy*) is a rule that a player uses to choose their next edge. In
99 the most general sense, a strategy can depend on the prefixes of paths. However, it is known
100 that for both discounted-sum and mean-payoff games with perfect information, memoryless and
101 deterministic optimal strategies always exist [Con92, EM79]. We refer to strategies that are both
102 memoryless and deterministic as *positional* strategies, to which we restrict our attention without
103 loss of generality. Formally, we define a strategy for the Min Player as a mapping $\pi_{\min} : V_{\min} \rightarrow V$
104 such that $(v, \pi_{\min}(v)) \in E$ for all $v \in V_{\min}$. Symmetrically, a strategy for the Max Player is a
105 mapping $\pi_{\max} : V_{\max} \rightarrow V$ such that $(v, \pi_{\max}(v)) \in E$ for all $v \in V_{\max}$. We denote the sets of all
106 valid positional strategies for the Min and Max Players as Π_{\min} and Π_{\max} , respectively. Given an

107 initial vertex $v \in V$ and a pair of positional strategies (π_{\min}, π_{\max}) , the game graph G produces
 108 a unique path, which we denote by $\rho(v, \pi_{\min}, \pi_{\max})$. Note that this path consists of a finite-path
 109 followed by a simple cycle infinitely repeated.

110 **Values.** It is well-known from the literature [Con92] that for any given γ and \mathbf{r} , every vertex
 111 $v \in V$ has a unique discounted-sum *value*, denoted by $\mu^{\gamma, \mathbf{r}}(v)$, which satisfies:

$$\mu^{\gamma, \mathbf{r}}(v) = \min_{\pi_{\min} \in \Pi_{\min}} \max_{\pi_{\max} \in \Pi_{\max}} u^{\gamma, \mathbf{r}}(\rho(v, \pi_{\min}, \pi_{\max})) = \max_{\pi_{\max} \in \Pi_{\max}} \min_{\pi_{\min} \in \Pi_{\min}} u^{\gamma, \mathbf{r}}(\rho(v, \pi_{\min}, \pi_{\max})).$$

112 Similarly, the mean-payoff game $G^{\text{mp}, \mathbf{r}}$ has a unique value for every $v \in V$, denoted by $\mu^{\text{mp}, \mathbf{r}}(v)$.

113 **Optimal strategies.** A strategy $\pi_{\min}^* \in \Pi_{\min}$ is an *optimal min strategy* for the discounted-sum
 114 game $G^{\gamma, \mathbf{r}}$ if it guarantees that the discounted-sum payoff never exceeds the discounted-sum value,
 115 regardless of the opponent's strategy; i.e., for all starting vertices $v \in V$ and all $\pi_{\max} \in \Pi_{\max}$, we
 116 have $u^{\gamma, \mathbf{r}}(\rho(v, \pi_{\min}^*, \pi_{\max})) \leq \mu^{\gamma, \mathbf{r}}(v)$. The notion of optimal max strategies is defined symmetrically.
 117 We say that $(\pi_{\min}^*, \pi_{\max}^*)$ is a pair of optimal strategies for $G^{\gamma, \mathbf{r}}$ if π_{\min}^* is an optimal min strategy
 118 and π_{\max}^* is an optimal max strategy for the discounted-sum objective. The notion of optimal
 119 strategies for mean-payoff games is defined analogously.

120 3 Smoothed model

121 The current best-known deterministic algorithms for discounted-sum and mean-payoff games have
 122 worst-case exponential complexity. However, in practice, instances of these games rarely exhibit
 123 worst-case behaviors. Hence the goal of this work is to study the problems through the lens of
 124 *smoothed analysis* [ST04], which bridges the gap between average-case and worst-case complexity
 125 by analyzing the performance of an algorithm under slight random perturbations of adversarial
 126 inputs. Throughout this section, we describe the smoothed model for discounted-sum games. The
 127 model for mean-payoff games is obtained similarly.

128 **Input formalization.** We follow the standard notion of smoothed model where the adversarial
 129 input is normalized. In discounted-sum and mean-payoff games, the rewards can be shifted and
 130 scaled to ensure that they lie in the interval $[-1, 1]$. In our smoothed analysis framework, an
 131 adversary first specifies the graph $G = (V_{\min}, V_{\max}, E)$ and the discount factor $\gamma \in (0, 1)$. The
 132 adversary also chooses an initial, worst-case reward vector $\mathbf{r}_0 \in [-1, 1]^E$. That is, the adversary
 133 can assign any initial reward to any edge, provided it is bounded within the interval $[-1, 1]$. After
 134 the adversary fixes G , γ , and \mathbf{r}_0 , a random continuous perturbation is applied to the rewards.
 135 Specifically, we define a noise vector $\xi \in \mathbb{R}^E$. For broad application, we do not restrict ξ to a
 136 specific distribution. Instead, we require the perturbation to satisfy a set of standard regularity
 137 properties, parameterized by a maximum density bound $\phi > 0$ and a tail decay parameter $\theta > 0$.

138 **Definition 2** (Smoothed model). The random noise vector $\xi \in \mathbb{R}^E$ satisfies the following properties:

- 139 1. *Independence:* The coordinates ξ_a for all $a \in E$ are mutually independent.
- 140 2. *Density bound ϕ :* For all $a \in E$, the probability distribution of ξ_a is absolutely continuous with
 141 respect to the Lebesgue measure. Furthermore, its probability density function is uniformly
 142 bounded by ϕ .
- 143 3. *Tail bound θ :* For all $a \in E$ and $t > 0$, the tail probability satisfies $\mathbb{P}[|\xi_a| > t] \leq \exp(-\theta t)$.

144 **Special cases:** Our general smoothed model captures many standard continuous distributions:

- 145 • *Gaussian noise:* For $\xi_a \sim \mathcal{N}(0, \sigma^2)$, the conditions hold with $\phi = \frac{1}{\sigma\sqrt{2\pi}}$ and $\theta = \frac{\sqrt{2/\pi}}{\sigma}$.
- 146 • *Laplacian noise:* For $\xi_a \sim \text{Laplace}(0, b)$, the conditions hold with $\phi = \frac{1}{2b}$ and $\theta = \frac{1}{b}$.
- 147 • *Uniform noise:* For $\xi_a \sim \text{Uniform}(-b, b)$, the conditions hold with $\phi = \frac{1}{2b}$ and $\theta = \frac{1}{b}$.

148 **Problem statement.** Given the input consisting of a game graph G , discount factor γ , and the
 149 normalized reward vector \mathbf{r}_0 , the algorithmic problem is to compute optimal strategies and values
 150 for the *smoothed game* $G^{\gamma, \mathbf{r}}$, where the final perturbed reward vector $\mathbf{r} \in \mathbb{R}^E$ is given by $\mathbf{r} = \mathbf{r}_0 + \xi$.

151 **Probability measure.** The random noise ξ induces a unique probability measure over the space
 152 of perturbed reward vectors \mathbf{r} . We define the probability measure as $\mathbb{P}_\xi[\cdot]$ and the corresponding
 153 expectation under this measure as $\mathbb{E}_\xi[\cdot]$. Since for any edge a the adversarial base rewards are
 154 bounded such that $\mathbf{r}_0(a) \in [-1, 1]$, the tail concentration property guarantees that the magnitude
 155 of the final perturbed reward satisfies

$$\mathbb{P}_\xi[|\mathbf{r}(a)| > 1 + t] \leq \exp(-\theta t). \quad (1)$$

156 **Polynomial-time smoothed complexity.** Finally, we define the smoothed complexity of an
 157 algorithm for this problem. Let $\mathcal{T}(G, \gamma, \mathbf{r}_0 + \xi)$ denote the total number of arithmetic operations
 158 of an algorithm on the perturbed game $G^{\gamma, \mathbf{r}}$. We evaluate the algorithm’s performance under two
 159 standard types of smoothed guarantees:

160 1. *High probability guarantee:* An algorithm achieves *smoothed polynomial time with high prob-*
 161 *ability* if, for any failure probability $\epsilon \in (0, 1)$, its runtime is bounded by a polynomial in
 162 the game dimensions (n and m), the noise parameters (ϕ and $1/\theta$), and the inverse failure
 163 probability ($1/\epsilon$), with probability at least $1 - \epsilon$. Formally, for any graph G , discount factor
 164 γ , initial reward vector $\mathbf{r}_0 \in [-1, 1]^E$:

$$\mathbb{P}_\xi \left[\mathcal{T}(G, \gamma, \mathbf{r}_0 + \xi) \leq \text{poly} \left(n, m, \phi, \frac{1}{\theta}, \frac{1}{\epsilon} \right) \right] \geq 1 - \epsilon.$$

165 2. *Expected runtime guarantee:* An algorithm achieves *expected smoothed polynomial time* if its
 166 expected runtime over the random noise is bounded by a polynomial in n, m, ϕ , and $1/\theta$.
 167 Formally, for any graph G , discount factor γ , and initial reward vector $\mathbf{r}_0 \in [-1, 1]^E$:

$$\mathbb{E}_\xi \left[\mathcal{T}(G, \gamma, \mathbf{r}_0 + \xi) \right] \leq \text{poly} \left(n, m, \phi, \frac{1}{\theta} \right).$$

168 *Remark 1* (High probability vs. expectation guarantee). We clarify the two types of guarantees
 169 above in the context of smoothed complexity. In the design of randomized algorithms, a high-
 170 probability bound is often considered a stronger guarantee because the algorithm can draw fresh
 171 randomness and restart upon failure. In smoothed analysis, however, the random perturbation is
 172 fixed exactly once, making restarts impossible. Consequently, an expected polynomial runtime is a
 173 strictly stronger mathematical guarantee. As specifically noted in [ST09], a high-probability poly-
 174 nomial bound follows from an expected polynomial bound via Markov’s inequality. The converse,
 175 however, does not hold: an algorithm may run efficiently on $1 - \epsilon$ fraction of the perturbations,
 176 yet its runtime could grow exponentially on the remaining ϵ fraction of bad instances, resulting in
 177 exponential expected runtime. Thus, bounding the expected runtime is the desired and stronger
 178 guarantee in the context of smoothed analysis. Thus the main goal of the study is to provide above
 179 guarantees, especially expected runtime guarantee, for discounted-sum and mean-payoff games.

180 4 Overview of Techniques

181 *Previous result and technique.* The previous result [LS24] in the literature for polynomial-time
 182 smoothed complexity for games with mean-payoff and discounted-sum objectives is restrictive in
 183 two ways: First, the class of game graphs is restricted to ergodic game graphs, where all vertices
 184 have the same value for mean-payoff objectives. Second, the result only provides smoothed high-
 185 probability guarantee and not the more general expected guarantee. The main algorithm for the
 186 smoothed analysis is the policy-iteration algorithm, which is simpler for ergodic mean-payoff games,

187 as compared to general mean-payoff games. While for general mean-payoff games, the policy-
 188 iteration algorithm for mean-payoff objectives relies on (i) values and (ii) potential or bias; for
 189 ergodic games, since all vertices have the same value, the only relevant concept is the bias. The
 190 results of [LS24] depend on the notion of bias-induced optimal strategies, which is only relevant
 191 for ergodic games. The key mathematical techniques in [LS24] along with bias-induced optimal
 192 strategies are (i) the notion of condition number, and (ii) geometric interpretation of reward vectors
 193 and hyper-plane separation induced by strategy profiles.

194 *Our techniques:* Our analysis provides the general result removing all the above restrictions. We
 195 present our key technical contributions below.

- 196 1. *Ingredient 1:* First, we provide a quantitative decomposition of discounted-sum values with
 197 respect to mean-payoff values (Lemma 2) and then present a Lipschitz continuity result with
 198 respect to switches of edges from optimal strategies (Lemma 3). These results hold for all
 199 game graphs, and are neither specific to the smoothed model nor a probabilistic guarantee.
- 200 2. *Ingredient 2:* Second, we present a probabilistic separation between (i) the value of a
 201 vertex and (ii) the value obtained by playing a sub-optimal edge and then playing opti-
 202 mally (Lemma 4); which we call the separation of the sub-optimality gap.
- 203 3. *Ingredient 3:* Finally, for general mean-payoff games, the policy iteration algorithm depends
 204 on both the values and the bias. Our final key ingredient is a mean-payoff value separation
 205 bound, which shows that for every vertex pair either the value is the same or there is a
 206 separation bound for the values (Lemma 5).

207 Ingredient 1 is the foundational result which is the basis of all results. We remark that with
 208 Ingredient 1 and Ingredient 2 we have an alternative proof for ergodic games with high-probability
 209 guarantee. However, for brevity, we focus on general game graphs for this paper. Ingredient 3
 210 along with the other two ingredients is the main idea for generalization from ergodic to all game
 211 graphs, for high-probability guarantee. For the extension from high-probability to expected runtime
 212 guarantee we present another technical contribution:

- 213 4. *Ingredient 4:* We consider the Blackwell-optimal discount factor γ_{bw} such that optimal strate-
 214 gies for discounted-sum are also optimal for the mean-payoff objectives. Let the inverse gap
 215 be $1/(1 - \gamma_{\text{bw}})$. The runtime complexity of the policy-iteration algorithm is dominated by the
 216 inverse gap, and in general the expected value of this inverse gap may be exponential. We
 217 therefore consider a truncated version of this inverse gap, which is the minimum of the inverse
 218 gap and the number of strategies, and show a polynomial upper bound for the expectation
 219 of the truncated version (Lemma 9).

220 Our ingredients tightly relate discounted-sum and mean-payoff games and hence we view mean-
 221 payoff games through the lens of discounted-sum games and Blackwell-optimality. Hence our main
 222 analysis focuses on discounted-sum games with arbitrary discount factors. We start with the formal
 223 background on Blackwell-optimality and policy iteration algorithm for discounted-sum games.

224 5 Background on Blackwell Optimality and Policy Iteration

225 In this section, we provide some useful theoretical and algorithmic tools from the literature. We
 226 first recall the concept of Blackwell optimality, which is a fundamental result in the theory of
 227 dynamic programming and games that establishes the existence of strategies, namely *Blackwell*
 228 *optimal strategies*, that are uniformly optimal for all discount factors sufficiently close to 1 [Bla62,
 229 BK76]. This property provides an important bridge between discounted-sum and mean-payoff
 230 games. We then describe the classic policy iteration algorithm for discounted-sum games to compute
 231 the optimal strategies and game values.

232 **Blackwell threshold $\gamma_{\text{bw}}(\mathbf{r})$ and Blackwell-optimal strategies.** Formally, given a specified

233 reward vector \mathbf{r} , a strategy $\pi_{\min}^{\text{bw}} \in \Pi_{\min}$ is a *Blackwell optimal min strategy* if there exists a threshold
 234 discount factor $\gamma_0 \in [0, 1)$ such that for all $\gamma \in [\gamma_0, 1)$, the strategy π_{\min}^{bw} is an optimal min strategy
 235 for the discounted-sum game $G^{\gamma, \mathbf{r}}$. Symmetrically, a strategy $\pi_{\max}^{\text{bw}} \in \Pi_{\max}$ is a *Blackwell optimal*
 236 *max strategy* with respect to \mathbf{r} if there exists a threshold discount factor $\gamma_0 \in [0, 1)$ such that for
 237 all $\gamma \in [\gamma_0, 1)$, π_{\max}^{bw} is an optimal max strategy for $G^{\gamma, \mathbf{r}}$. We define the *Blackwell threshold* with
 238 respect to \mathbf{r} , denoted by $\gamma_{\text{bw}}(\mathbf{r}) \in [0, 1)$, as the infimum over all such valid thresholds.

239 **Asymptotic relationship between discounted-sum and mean-payoff values.** By definition,
 240 for any $\gamma \in [\gamma_{\text{bw}}(\mathbf{r}), 1)$, a pair of Blackwell optimal strategies $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ is guaranteed to be opti-
 241 mal for the discounted-sum game $G^{\gamma, \mathbf{r}}$. Furthermore, these strategies are of particular theoretical
 242 importance because they simultaneously act as optimal strategies for the corresponding mean-payoff
 243 game $G^{\text{mp}, \mathbf{r}}$, providing a crucial bridge between discounted-sum and mean-payoff games [BK76].
 244 Formally, for any reward vector $\mathbf{r} \in \mathbb{R}^E$ and vertex $v \in V$, the normalized discounted-sum value
 245 converges exactly to the mean-payoff value, i.e.,

$$\mu^{\text{mp}, \mathbf{r}}(v) = \lim_{\gamma \uparrow 1} (1 - \gamma) \mu^{\gamma, \mathbf{r}}(v). \quad (2)$$

246 **Policy iteration.** To compute the optimal strategies and values in the discounted-sum game $G^{\gamma, \mathbf{r}}$,
 247 we describe the classic policy iteration algorithm, formally presented in Algorithm 1 in Appendix B.
 248 The algorithm operates from the perspective of one of the players; without loss of generality, we
 249 formulate it for the Min Player. Starting from an arbitrary initial strategy, the algorithm iteratively
 250 refines the strategy π_{\min} through alternating phases of policy evaluation and policy improvement.

251 • *Policy evaluation phase:* The algorithm fixes the Min Player’s current strategy π_{\min} . With
 252 π_{\min} fixed, the game reduces to a graph with only choices for the Max Player (aka deterministic
 253 MDP). Graphs (or deterministic MDPs) with discounted-sum objectives can be solved in
 254 strongly polynomial time [PY15]. The resulting vector μ represents the expected discounted-
 255 sum payoff under the current strategies.

256 • *Policy improvement phase:* The algorithm searches for *improving switches*—local edge
 257 changes that strictly decrease the current evaluated payoff for the Min Player. If no improv-
 258 ing switch exists for any vertex in V_{\min} , the algorithm terminates. Otherwise, the algorithm
 259 updates π_{\min} by simultaneously applying the best improving switch at every possible vertex.
 260 This specific greedy update mechanism is known as Howard’s all-switch rule.

261 In other words, the policy-iteration algorithm iterates over strategies of a player by greedily switch-
 262 ing over local improvements.

263 **Correctness and complexity.** The key correctness proof establishes that when no local-
 264 improvement switches are possible, then the obtained strategy is globally optimal [Con92]. Since
 265 the total number of strategies is at most n^n and each update strictly decreases the value vector, the
 266 algorithm is guaranteed to converge in finite time [FV96]. Specifically, for discounted-sum games,
 267 the upper bound on the number of iterations depends on the discount factor γ . The worst-case
 268 strongly polynomial bound for a constant discount factor from the literature, which is the basic
 269 result we use in subsequent analysis, is as follows:

270 **Proposition 1** (HMZ13, Theorem 7.5). *Let $\gamma \in [0, 1)$. For any game $G^{\gamma, \mathbf{r}}$ with n vertices and m
 271 edges, the number of iterations that PI algorithm (Algorithm 1) requires is at most $6 \frac{m}{1-\gamma} \ln\left(\frac{n}{1-\gamma}\right)$.*

272 Each iteration of PI performs a single policy evaluation, which requires at most $\mathcal{O}(n^3 m^2 \log^2 n)$
 273 arithmetic operations [PY15]. While this algorithm performs well in practice, the worst-case bound
 274 of Proposition 1 is exponential as γ approaches 1 (e.g., $\gamma = 1 - 2^{-n}$).

6 Mathematical Ingredients 1-3

In this section, we present the details related to the first three ingredients mentioned in Section 4; full proofs are provided in Appendix C.

6.1 Ingredient 1: Value Decomposition and Lipschitz Continuity

In this subsection, we present conceptual ideas behind Ingredient 1 from Section 4.

- We first establish a value decomposition for discounted-sum games under Blackwell optimality (Lemma 2), which provides a useful connection between discounted-sum and mean-payoff games. As compared to Equation (2) which only provides convergence-in-the-limit guarantee, this result provides a quantitative bound on the difference of the discounted-sum and mean-payoff values. While a general characterization of discounted-sum and mean-payoff value is provided in [BK76] for concurrent games (with simultaneous interaction), we present a quantitative characterization albeit for the simpler class of turn-based games.
- We then define the suboptimality gap, which quantifies the suboptimality incurred by switching a single edge in the optimal strategy. We show that whenever the two endpoints of an edge share the same mean-payoff value, the suboptimality gap is Lipschitz continuous with respect to the discount factor γ over the Blackwell threshold (Lemma 3)

Lemma 2 (Value decomposition under Blackwell optimality). *Consider a game graph G with reward vector $\mathbf{r} \in \mathbb{R}^E$. Let $v \in V$ be a vertex.*

1. Decomposition: *For all $\gamma \in (\gamma_{bw}(\mathbf{r}), 1)$, the discounted-sum value admits the following decomposition, where $w_v^{\mathbf{r}}(\gamma)$ is a rational function of γ :*

$$\mu^{\gamma, \mathbf{r}}(v) = \frac{\mu^{\text{mp}, \mathbf{r}}(v)}{1 - \gamma} + w_v^{\mathbf{r}}(\gamma).$$

2. Magnitude and derivative: *For all $\gamma \in (\gamma_{bw}(\mathbf{r}), 1)$, the magnitude and the derivative of $w_v^{\mathbf{r}}(\gamma)$ are bounded by:*

$$|w_v^{\mathbf{r}}(\gamma)| \leq 4n \|\mathbf{r}\|_{\infty} \quad \text{and} \quad \left| \left(\frac{d}{d\gamma} w_v^{\mathbf{r}}(\gamma) \right) \right| \leq 7n^2 \|\mathbf{r}\|_{\infty}.$$

Proof sketch. Let $\mu \triangleq \mu^{\text{mp}, \mathbf{r}}(v)$. Fix a Blackwell optimal strategy pair, under which the path from v consists of a transient path of length $t \leq n$ followed by a cycle of length $c \leq n$. Let r_k denote the reward at step k . Then μ is the average cycle reward, so $\sum_{j=0}^{c-1} (r_{t+j} - \mu) = 0$. Writing each reward as $(r_k - \mu) + \mu$ in $\mu^{\gamma, \mathbf{r}}(v) = \sum_{j=0}^{t-1} r_j \gamma^j + \frac{\gamma^t}{1 - \gamma^c} \sum_{j=0}^{c-1} r_{t+j} \gamma^j$ shows that the μ terms contribute $\mu/(1 - \gamma)$. The cycle cancellation removes this singular factor from the remaining terms, which define the rational function $w_v^{\mathbf{r}}(\gamma)$. For magnitude and derivative, $w_v^{\mathbf{r}}(\gamma)$ contains at most $2n$ deviations, each bounded by $2\|\mathbf{r}\|_{\infty}$ and multiplied by a coefficient in $[0, 1]$, giving $|w_v^{\mathbf{r}}(\gamma)| \leq 4n\|\mathbf{r}\|_{\infty}$. Differentiating these coefficients bounds the transient contribution by $n^2\|\mathbf{r}\|_{\infty}$ and the cycle contribution by $6n^2\|\mathbf{r}\|_{\infty}$, yielding the derivative bound. \square

Suboptimality gap $\Delta^{\gamma, \mathbf{r}}$. We define the *suboptimality gap* for the discounted-sum game $G^{\gamma, \mathbf{r}}$ as follows. For all edges $(v, v') \in E$,

$$\Delta^{\gamma, \mathbf{r}}(v, v') = \mathbf{r}(v, v') + \gamma \mu^{\gamma, \mathbf{r}}(v') - \mu^{\gamma, \mathbf{r}}(v);$$

i.e., the quantity $\mathbf{r}(v, v') + \gamma \mu^{\gamma, \mathbf{r}}(v')$ is the value for switching via the edge, and $\mu^{\gamma, \mathbf{r}}(v)$ is the optimal value; and hence the function represents the relative value for switching via the edge.

Main property. For any reward vector $\mathbf{r} \in \mathbb{R}^E$ such that $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$, we show that the suboptimality gap is Lipschitz continuous in γ over the Blackwell threshold. In contrast, the value vector $\mu^{\gamma, \mathbf{r}}(v)$ may not be Lipschitz continuous in γ because the singular term $\frac{\mu^{\text{mp}, \mathbf{r}}(v)}{1 - \gamma}$ in the value decomposition (Lemma 2) diverges as γ approaches 1.

314 **Lemma 3** (Lipschitz continuity). *Let $\mathbf{r} \in \mathbb{R}^E$ be any reward vector. For any edge $a = (v, v') \in E$*
 315 *where $\mu^{\text{mp},\mathbf{r}}(v) = \mu^{\text{mp},\mathbf{r}}(v')$, the map $\gamma \mapsto \Delta^{\gamma,\mathbf{r}}(v, v')$ is $18n^2\|\mathbf{r}\|_\infty$ -Lipschitz continuous over the*
 316 *interval $(\gamma_{bw}(\mathbf{r}), 1)$.*

317 *Proof sketch.* Substituting the value decomposition from [Lemma 2](#) into the suboptimality gap,
 318 the shared mean-payoff value cancels the singular terms, giving $\Delta^{\gamma,\mathbf{r}}(v, v') = \mathbf{r}(v, v') - \mu^{\text{mp},\mathbf{r}}(v) +$
 319 $\gamma w_{v'}^{\mathbf{r}}(\gamma) - w_v^{\mathbf{r}}(\gamma)$. Therefore, using the magnitude and derivative bounds from [Lemma 2](#), we have
 320 $\left| \frac{d}{d\gamma} \Delta^{\gamma,\mathbf{r}}(v, v') \right| \leq 4n\|\mathbf{r}\|_\infty + 14n^2\|\mathbf{r}\|_\infty \leq 18n^2\|\mathbf{r}\|_\infty$, which proves the result. \square

321 **Gap at one $\Delta^{1,\mathbf{r}}$.** The Lipschitz continuity result implies that we can extend the definition of
 322 $\Delta^{\gamma,\mathbf{r}}$ to $\gamma = 1$ as follows. For any edge $(v, v') \in E$ where $\mu^{\text{mp},\mathbf{r}}(v) = \mu^{\text{mp},\mathbf{r}}(v')$, we have

$$\Delta^{1,\mathbf{r}}(v, v') \triangleq \lim_{\gamma \uparrow 1} \Delta^{\gamma,\mathbf{r}}(v, v').$$

323 Following from [Lemma 3](#), the limit exists and is finite. We also refer to $\Delta^{1,\mathbf{r}}$ as the *suboptimality*
 324 *gap* of the mean-payoff game $G^{\text{mp},\mathbf{r}}$.

325 6.2 Ingredient 2: Suboptimality Gap Separation

326 We now establish a probabilistic separation result for the suboptimality gap under the smoothed
 327 model (Ingredient 2 from [Section 4](#)). For any fixed edge $a = (v, v')$, the event that (i) vertices v and
 328 v' have the same mean-payoff value and (ii) the suboptimality gap at $\gamma = 1$ is small and non-zero,
 329 has low probability.

330 **Lemma 4** (Suboptimality separation). *Let G be any game graph and let $a = (v, v') \in E$ be any*
 331 *fixed edge. Under the smoothed model, for any $\delta > 0$, we have*

$$\mathbb{P}_\xi \left[\mu^{\text{mp},\mathbf{r}}(v) = \mu^{\text{mp},\mathbf{r}}(v') \text{ and } 0 < |\Delta^{1,\mathbf{r}}(a)| < \delta \right] \leq 2\delta\phi.$$

332 *Proof sketch.* Consider the event $\mathcal{E} \triangleq \{ \mathbf{r} \in \mathbb{R}^E : \mu^{\text{mp},\mathbf{r}}(v) = \mu^{\text{mp},\mathbf{r}}(v') \text{ and } 0 < |\Delta^{1,\mathbf{r}}(a)| < \delta \}$. If a
 333 is the only outgoing edge from v , then $\Delta^{1,\mathbf{r}}(a) = 0$ whenever well-defined, so \mathcal{E} is empty. Otherwise,
 334 fix the rewards \mathbf{r}_{-a} on all edges except a and write $x \triangleq \mathbf{r}(a)$. For any x such that $(x; \mathbf{r}_{-a}) \in \mathcal{E}$, the
 335 edge a is strictly suboptimal for all γ sufficiently close to 1. Deleting a therefore leaves the value
 336 vector unchanged and preserves the shared mean-payoff value of v and v' . Thus,

$$\Delta^{1,\mathbf{r}}(a) = x + \lim_{\gamma \uparrow 1} (\gamma \mu^{\gamma,\mathbf{r}-a}(v') - \mu^{\gamma,\mathbf{r}-a}(v)).$$

337 By [Lemma 2](#), the limit term is finite and independent of x . Hence the set of x values producing \mathcal{E}
 338 has length at most 2δ , and therefore conditional probability at most $2\delta\phi$. Since this bound holds
 339 for every fixed \mathbf{r}_{-a} , the unconditional bound follows. \square

340 6.3 Ingredient 3: Mean-payoff Value Separation

341 We establish Ingredient 3 from [Section 4](#): under the smoothed model, the mean-payoff values of any
 342 two vertices are either the same or separated by an inverse-polynomial gap with high probability.

343 **Lemma 5** (Mean-payoff separation). *Let G be any game graph and let $v, v' \in V$. Under the*
 344 *smoothed model, for any $B > 0$, we have*

$$\mathbb{P}_\xi \left[0 < |\mu^{\text{mp},\mathbf{r}}(v) - \mu^{\text{mp},\mathbf{r}}(v')| < B \right] \leq 2Bm\phi.$$

345 *Proof sketch.* With probability 1 the Blackwell optimal pair is unique; fix it and let $C_v^{\mathbf{r}}$ be the
 346 recurrent cycle reached from v , of average reward $\mu^{\text{mp},\mathbf{r}}(v)$. If $\mu^{\text{mp},\mathbf{r}}(v) > \mu^{\text{mp},\mathbf{r}}(v')$, then no edge of
 347 $C_v^{\mathbf{r}}$ is reachable from v' under the optimal Min strategy π_{\min}^{bw} , since otherwise Max could reach and
 348 repeat $C_v^{\mathbf{r}}$ from v' and secure $\mu^{\text{mp},\mathbf{r}}(v) > \mu^{\text{mp},\mathbf{r}}(v')$ against π_{\min}^{bw} . Assume $0 < \mu^{\text{mp},\mathbf{r}}(v) - \mu^{\text{mp},\mathbf{r}}(v') <$

349 B , fix an edge $a \in E(C_v^{\mathbf{r}})$, condition on \mathbf{r}_{-a} , and write $x \triangleq \mathbf{r}(a)$. Since a is unreachable from v' , the
350 value $\mu^{\text{mp},(x;\mathbf{r}-a)}(v')$ is a constant α independent of x , whereas $x \mapsto \mu^{\text{mp},(x;\mathbf{r}-a)}(v)$ is nondecreasing
351 with derivative at least $1/n$ almost everywhere on the set where a lies on the cycle reached from v .
352 This set therefore maps into $(\alpha, \alpha + B)$ and has Lebesgue measure at most nB ; multiplying by the
353 density bound ϕ , summing over the m edges, and adding the symmetric case gives $2mnB\phi$. \square

354 7 Discounted-sum games

355 In this section, we focus our attention on discounted-sum games under the smoothed model. To
356 establish smoothed polynomial runtime, we first derive a tail bound on the Blackwell threshold
357 $\gamma_{\text{bw}}(\mathbf{r})$ by combining the mathematical ingredients from Section 6 (Section 7.1). We then introduce
358 a restarting variant of policy iteration and establish a total iteration bound for it in terms of $\gamma_{\text{bw}}(\mathbf{r})$
359 (Section 7.2). Combining this bound with the tail bound yields smoothed polynomial runtime with
360 high probability (Section 7.3). Combining it with a bound on the expected truncated inverse gap
361 yields smoothed polynomial expected runtime (Section 7.4). The full proofs of the results in this
362 section are provided in Appendix D.

363 7.1 Tail bound

364 In this subsection, we establish a tail bound on the Blackwell threshold $\gamma_{\text{bw}}(\mathbf{r})$ under the smoothed
365 model. The resulting bound shows that $\gamma_{\text{bw}}(\mathbf{r})$ is bounded away from 1 with high probability and
366 serves as the foundation for both the high probability and expected runtime analyses in the sequel.
367 We first define the variable K which is a polynomial in the input parameters m , n , and ϕ :

$$K \triangleq 92m^2n^2(\phi + 1). \quad (3)$$

368

369 **Lemma 6** (Tail bound). *Let G be any game graph. Under the smoothed model, for any $x \geq 1$,*
370 *with the bounding function defined as $M(x) \triangleq 1 + \frac{1}{\phi} \ln(mx)$, we have*

$$\mathbb{P}_{\xi}[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] \leq \frac{KM(x)}{x}.$$

371 *Proof sketch.* We combine the reward tail bound Eq. (1) with the suboptimality gap and mean-
372 payoff-value separation events from Lemma 4 and Lemma 5. Condition on the intersection of
373 these events, Lemma 2 and Lemma 3 imply that every non-optimal edge retains its sign for all
374 $\gamma > 1 - 1/x$, and hence $\gamma_{\text{bw}}(\mathbf{r}) \leq 1 - 1/x$. Choosing the separation parameters proportional to
375 $M(x)/x$ and applying union bounds over the edges shows that the complement of this event has
376 probability at most $KM(x)/x$. \square

377 7.2 Restarting Policy Iteration

378 In this subsection, we introduce the restarting variant of policy iteration called RePI, presented
379 in full detail in Algorithm 2 in Appendix B. We bound its total iteration count in Lemma 7:
380 deterministically by the number of positional strategies for every reward vector, and, under the
381 smoothed model, by an inverse-polynomial factor of the Blackwell gap $1 - \gamma_{\text{bw}}(\mathbf{r})$ with probability 1.

382 **Informal description.** Our algorithm sweeps through a geometric sequence of thresholds $\gamma_k =$
383 $1 - e^{-k}$ for $k = 0, 1, \dots, \lceil n \ln n \rceil$. For each k in turn, if the target discount factor γ already satisfies
384 $\gamma \leq \gamma_k$, the algorithm directly executes the standard PI procedure on γ . Otherwise, it solves the
385 game with the discount factor γ_k to obtain a candidate strategy pair, computes the payoff of this
386 pair under the target discount factor γ , and checks for any unilateral improving switches. If none
387 exists, the candidate is optimal for the discount factor γ and the algorithm terminates. If the
388 test fails, the algorithm proceeds to the next threshold γ_{k+1} and retries. If every threshold in the
389 sequence fails, then the algorithm falls back to running the standard PI procedure on γ .

390 **Lemma 7** (Total iteration bound). *Let G be any game graph with $n \geq 2$, let $\gamma \in [0, 1)$, and let*
 391 *$\mathbf{r} \in \mathbb{R}^E$. Write $N_{iter}(\mathbf{r})$ for the total number of iterations executed by $\text{RePI}(G, \gamma, \mathbf{r})$ (Algorithm 2),*
 392 *summed over all calls to PI . Then $N_{iter}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2)n^n$ for every \mathbf{r} , and, under the smoothed*
 393 *model, with probability 1,*

$$N_{iter}(\mathbf{r}) \leq 27 \cdot \frac{m}{1 - \gamma_{bw}(\mathbf{r})} \ln\left(\frac{e \cdot n}{1 - \gamma_{bw}(\mathbf{r})}\right).$$

394 7.3 High Probability Guarantees

395 In this subsection, we combine the total iteration bound of Lemma 7 with the tail bound from
 396 Section 7.1 to show that, with probability at least $1 - \epsilon$, RePI terminates after a number of iterations
 397 polynomial in the game parameters, the density bound, the inverse of tail bound, and the inverse
 398 failure probability $1/\epsilon$.

399 **Theorem 8.** *The restarting PI algorithm achieves smoothed polynomial time with high probability*
 400 *guarantee for discounted-sum games.*

401 *Proof sketch.* Fix a failure probability $\epsilon \in (0, 1)$. The tail bound in Lemma 6 shows that the inverse
 402 Blackwell gap $1/(1 - \gamma_{bw}(\mathbf{r}))$ is bounded by a polynomial in $n, m, \phi, 1/\theta$, and $1/\epsilon$, with probability
 403 at least $1 - \epsilon$. Conditioned on this event, substituting this bound into the total iteration bound of
 404 Lemma 7 yields a polynomial bound on the total number of iterations executed by RePI . The full
 405 proof is given in Lemma 16 in Appendix D. \square

406 7.4 Expected Runtime Guarantees

407 To bound the expected number of iterations of RePI , we first define the notion of truncated inverse
 408 gap and show that its expectation is bounded by a polynomial in the game parameters, the density
 409 bound, and the inverse of the tail bound in Lemma 9 (Ingredient 4 from Section 4). We then
 410 combine this bound with the total iteration bound of Lemma 7 to obtain the desired polynomial
 411 bound on the expected number of iterations.

412 **Truncated inverse gap.** For any reward vector $\mathbf{r} \in \mathbb{R}^E$, define

$$Y(\mathbf{r}) \triangleq \min(1/(1 - \gamma_{bw}(\mathbf{r})) \ln(1/(1 - \gamma_{bw}(\mathbf{r}))), n^n).$$

413 **Lemma 9** (Bounding expected truncated inverse gap). *Let G be any game graph with $n \geq 2$. Let*
 414 *K be defined in Eq. (3). Under the smoothed model, we have*

$$\mathbb{E}_\xi[Y(\mathbf{r})] \leq \frac{6(\theta + 1)}{\theta} K(n \ln m)^3.$$

415 *Proof sketch.* By Lemma 6, the inverse Blackwell gap has tail probability of $\mathcal{O}(\ln(x)/x)$. Integrating
 416 this tail up to the truncation point n^n yields the stated polynomial bound. \square

417 **Theorem 10.** *The restarting PI algorithm achieves smoothed polynomial time with expected run-*
 418 *time guarantee for discounted-sum games.*

419 *Proof sketch.* By Lemma 7, the runtime of RePI is bounded by the inverse Blackwell gap (up to
 420 polynomial factors of n and m), while its worst-case fallback is bounded by the finite number of
 421 strategies. These two bounds are captured by the truncated inverse gap $Y(\mathbf{r})$. Therefore, the
 422 polynomial expectation of $Y(\mathbf{r})$ from Lemma 9 yields a polynomial expected runtime. The full
 423 proof is given in Lemma 17 in Appendix D. \square

424 As established in Section 5, any Blackwell optimal strategy for a discounted-sum game is an
 425 optimal strategy for the corresponding mean-payoff game [BK76]. Therefore, our smoothed analysis
 426 for computing Blackwell optimal strategies directly yields a smoothed polynomial-time algorithm
 427 for mean-payoff games, presented in Corollary 11. The full proof is given in Appendix E.

428 **Corollary 11.** *The restarting mean-payoff PI algorithm RePI_{mp} achieves smoothed polynomial time*
429 *for mean-payoff games, with both high-probability and expected runtime guarantees.*

430 Acknowledgement

431 The authors are thankful for the early discussion with Raimundo Saona and Elahe Tohidi. This
432 paper receives funding from ERC CoG 863818 (ForM-SMArt) and Austrian Science Fund (FWF)
433 10.55776/COE12.

434 References

- 435 [ACSSU24] Ali Asadi, Krishnendu Chatterjee, Jakub Svoboda, and Raimundo Saona Urmeneta. Deterministic
436 sub-exponential algorithm for discounted-sum games with unary weights. In *Proceedings of*
437 *the 39th Annual ACM/IEEE Symposium on Logic in Computer Science*, pages 6:1–6:12, 2024.
438 [1](#)
- 439 [BCHJ09] Roderick Bloem, Krishnendu Chatterjee, Thomas A Henzinger, and Barbara Jobstmann. Better
440 quality in synthesis through quantitative objectives. In *International Conference on Computer*
441 *Aided Verification*, pages 140–156. Springer, 2009. [1](#)
- 442 [BEF⁺11] Endre Boros, Khaled Elbassioni, Mahmoud Fouz, Vladimir Gurvich, Kazuhisa Makino, and
443 Bodo Manthey. Stochastic mean payoff games: Smoothed analysis and approximation schemes.
444 In *International Colloquium on Automata, Languages, and Programming*, pages 147–158.
445 Springer, 2011. [1](#)
- 446 [BK76] Truman Bewley and Elon Kohlberg. The asymptotic theory of stochastic games. *Mathematics*
447 *of Operations Research*, 1(3):197–208, 1976. [5](#), [6.1](#), [7.4](#), [C.1](#), [C.1](#), [C.3](#), [E.1](#), [E.1](#)
- 448 [Bla62] David Blackwell. Discrete dynamic programming. *The Annals of Mathematical Statistics*,
449 33(2):719–726, 1962. [5](#), [C.1](#)
- 450 [BV07] Henrik Björklund and Sergei Vorobyov. A combinatorial strongly subexponential strategy im-
451 provement algorithm for mean payoff games. *Discrete Applied Mathematics*, 155(2):210–229,
452 2007. [1](#), [E.1](#), [9](#), [E.1](#)
- 453 [CIJ14] Krishnendu Chatterjee and Rasmus Ibsen-Jensen. The complexity of ergodic mean-payoff
454 games. In *International Colloquium on Automata, Languages, and Programming*, pages 122–133.
455 Springer, 2014. [1](#)
- 456 [CKS81] Ashok K Chandra, Dexter C Kozen, and Larry J Stockmeyer. Alternation. *Journal of the ACM*
457 *(JACM)*, 28(1):114–133, 1981. [1](#)
- 458 [CLSZ25] Krishnendu Chatterjee, David Lurie, Raimundo Saona, and Bruno Ziliotto. Uniform value and
459 decidability in ergodic blind stochastic games. *Mathematics of Operations Research*, 2025. in
460 press. [1](#)
- 461 [Con92] Anne Condon. The complexity of stochastic games. *Information and Computation*, 96(2):203–
462 224, 1992. [1](#), [2](#), [2](#), [5](#)
- 463 [CPV15] Krishnendu Chatterjee, Andreas Pavlogiannis, and Yaron Velner. Quantitative interprocedural
464 analysis. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles*
465 *of Programming Languages*, pages 539–551, 2015. [1](#)

- 466 [CY23] Miranda Christ and Mihalis Yannakakis. The smoothed complexity of policy iteration for markov
467 decision processes. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*,
468 pages 1890–1903, 2023. [1](#)
- 469 [DAH03] Luca De Alfaro, Thomas A Henzinger, and Rupak Majumdar. Discounting the future in systems
470 theory. In *International Colloquium on Automata, Languages, and Programming*, pages 1022–
471 1037. Springer, 2003. [1](#)
- 472 [DKZ19] Dani Dorfman, Haim Kaplan, and Uri Zwick. A faster deterministic exponential time algorithm
473 for energy games and mean payoff games. In *46th International Colloquium on Automata,
474 Languages, and Programming (ICALP 2019)*, pages 114:1–114:14. Schloss Dagstuhl–Leibniz-
475 Zentrum für Informatik, 2019. [1](#)
- 476 [EM79] Andrzej Ehrenfeucht and Jan Mycielski. Positional strategies for mean payoff games. *Internat-*
477 *ional Journal of Game Theory*, 8(2):109–113, 1979. [1](#), [2](#)
- 478 [FV96] Jerzy Filar and Koos Vrieze. *Competitive Markov decision processes*. Springer, 1996. [1](#), [5](#)
- 479 [Gil57] Dean Gillette. Stochastic games with zero stop probabilities. *Contributions to the Theory of*
480 *Games III*, 39:179–187, 1957. [1](#)
- 481 [GKK88] Vladimir A Gurvich, Alexander V Karzanov, and Leonid G Khachiyan. Cyclic games and an
482 algorithm to find minimax cycle means in directed graphs. *USSR Computational Mathematics*
483 *and Mathematical Physics*, 28(5):85–91, 1988. [1](#)
- 484 [HK66] Alan J Hoffman and Richard M Karp. On nonterminating stochastic games. *Management*
485 *Science*, 12(5):359–370, 1966. [1](#)
- 486 [HMZ13] Thomas Dueholm Hansen, Peter Bro Miltersen, and Uri Zwick. Strategy iteration is strongly
487 polynomial for 2-player turn-based stochastic games with a constant discount factor. *Journal of*
488 *the ACM (JACM)*, 60(1):1–16, 2013. [1](#)
- 489 [HZ98] Eric A Hansen and Shlomo Zilberstein. Heuristic search in cyclic and/or graphs. In *AAAI/IAAI*,
490 pages 412–418, 1998. [1](#)
- 491 [Koz21] Alexander Kozachinskiy. Polyhedral value iteration for discounted games and energy games. In
492 *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 600–
493 616. SIAM, 2021. [1](#)
- 494 [LS24] Bruno Loff and Mateusz Skomra. Smoothed analysis of deterministic discounted and mean-
495 payoff games. In *51st International Colloquium on Automata, Languages, and Programming*
496 *(ICALP 2024)*, pages 147:1–147:16. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2024. [1](#),
497 [4](#), [B](#)
- 498 [Lud95] Walter Ludwig. A subexponential randomized algorithm for the simple stochastic game problem.
499 *Information and computation*, 117(1):151–155, 1995. [1](#)
- 500 [MB85] Ambuj Mahanti and Amitava Bagchi. And/or graph heuristic search methods. *Journal of the*
501 *ACM (JACM)*, 32(1):28–51, 1985. [1](#)
- 502 [MHC03] Omid Madani, Steve Hanks, and Anne Condon. On the undecidability of probabilistic planning
503 and related stochastic optimization problems. *Artificial Intelligence*, 147(1-2):5–34, 2003. [1](#)
- 504 [PR89] Amir Pnueli and Roni Rosner. On the synthesis of a reactive module. In *Proceedings of the 16th*
505 *ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 179–190,
506 1989. [1](#)
- 507 [Put94] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
508 Wiley & Sons, 1994. [1](#)

- 509 [PY15] Ian Post and Yinyu Ye. The simplex method is strongly polynomial for deterministic markov
510 decision processes. *Mathematics of Operations Research*, 40(4):859–868, 2015. [5](#), [5](#)
- 511 [RW87] Peter J Ramadge and W Murray Wonham. Supervisory control of a class of discrete event
512 processes. *SIAM journal on control and optimization*, 25(1):206–230, 1987. [1](#)
- 513 [Sha53] Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–
514 1100, 1953. [1](#)
- 515 [ST04] Daniel A Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex
516 algorithm usually takes polynomial time. *Journal of the ACM (JACM)*, 51(3):385–463, 2004. [1](#),
517 [3](#)
- 518 [ST09] Daniel A Spielman and Shang-Hua Teng. Smoothed analysis: an attempt to explain the behavior
519 of algorithms in practice. *Communications of the ACM*, 52(10):76–84, 2009. [1](#)
- 520 [ZP96] Uri Zwick and Mike Paterson. The complexity of mean payoff games on graphs. *Theoretical*
521 *Computer Science*, 158(1-2):343–359, 1996. [1](#)

522 A Notation summary

523 **Table 1** collects the recurring notation of the paper. Symbols are grouped by role; entries point to
 524 the section or equation where they are first introduced.

<i>Game graph and games</i>	
$G = (V_{\min}, V_{\max}, E)$	game graph; $V = V_{\min} \uplus V_{\max}$, $n = V $, $m = E $
$G^{\gamma, \mathbf{r}}$	discounted-sum game on G with discount γ and rewards \mathbf{r}
$G^{\text{mp}, \mathbf{r}}$	mean-payoff game on G with rewards \mathbf{r}
<i>Strategies</i>	
π_{\min}, π_{\max}	positional strategies of the Min / Max player
Π_{\min}, Π_{\max}	sets of all positional strategies
$\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}}$	Blackwell optimal strategies (Section 5)
I_{\min}, I_{\max}	vertices with improving switches for Min / Max
<i>Rewards, perturbation, values</i>	
$\mathbf{r}_0 \in [-1, 1]^E$	adversarial (unperturbed) reward vector
$\xi \in \mathbb{R}^E$	noise vector; density bounded by ϕ , tail parameter θ
$\mathbf{r} = \mathbf{r}_0 + \xi \in \mathbb{R}^E$	perturbed reward vector (Section 3)
$\gamma \in [0, 1)$	discount factor
$\mu^{\gamma, \mathbf{r}}(v), \mu^{\text{mp}, \mathbf{r}}(v)$	discounted-sum / mean-payoff value at vertex v
$\gamma_{\text{bw}}(\mathbf{r})$	Blackwell threshold (Section 5)
<i>Algorithms</i>	
$\text{PI}(G, \gamma, \mathbf{r})$	policy iteration on $G^{\gamma, \mathbf{r}}$ (Section 5)
RePI	restarting policy iteration for discounted-sum games (Algorithm 2)
RePI_{mp}	restarting policy iteration for mean-payoff games (Algorithm 3)
$\gamma_k = 1 - e^{-k}$	geometric threshold schedule used by RePI / RePI_{mp}
<i>Analysis quantities</i>	
$\Delta^{\gamma, \mathbf{r}}(v, v')$	suboptimality gap of the edge (v, v')
$K \triangleq 92m^2n^2(\phi + 1)$	coefficient in the tail bound (Eq. (3))
$M(x) \triangleq 1 + \frac{1}{\theta} \ln(mx)$	tail bound function (Lemma 6)
$Y(\mathbf{r})$	truncated inverse gap (Section 7.4)
$\epsilon \in (0, 1)$	failure probability for high-probability guarantees

Table 1: Recurring notation used throughout the paper.

525 B Algorithm Pseudocodes

526 In this section, we present the pseudocodes of the policy iteration algorithms for discounted-sum
 527 games discussed in the main body: the classic policy iteration **PI** ([Algorithm 1](#)) and the restart-
 528 ing variant **RePI** ([Algorithm 2](#)). The second algorithm, [Algorithm 2](#), is related in spirit to the
 529 increasing-discount algorithm **IncreasingDiscountPI2** of [[LS24](#)], but the crucial difference is that [Al-](#)
 530 [gorithm 2](#) has an explicit fallback after a fixed finite sequence of thresholds. This gives a uniform
 531 deterministic bound on the running time, which is essential for our expected-runtime analysis.
 532 Thus, although **IncreasingDiscountPI2** yields a high-probability smoothed guarantee in [[LS24](#)], the

533 algorithm does not achieve a polynomial expected-running-time guarantee based on our techniques.
534 Moreover, IncreasingDiscountPI2 warm-starts policy iteration from the policies computed at the pre-
535 vious discount factor, whereas Algorithm 2 restarts policy iteration from scratch at each threshold;
536 this difference, however, is not used in our analysis.

Algorithm 1 Policy Iteration: $\text{PI}(G, \gamma, \mathbf{r})$

```

1: Input: Graph  $G = (V_{\min}, V_{\max}, E)$ , discount factor  $\gamma \in [0, 1)$ , reward vector  $\mathbf{r} \in \mathbb{R}^E$ .
2: Output: Optimal strategy pair  $(\pi_{\min}^*, \pi_{\max}^*)$  and value vector  $\mu^{\gamma, \mathbf{r}}$  for the game  $G^{\gamma, \mathbf{r}}$ .
3: Initialize an arbitrary positional strategy  $\pi_{\min} \in \Pi_{\min}$  for the Min Player.
4: loop
5:   // Policy Evaluation Phase
6:   Compute the value vector  $\mu \in \mathbb{R}^V$  and Max Player's best-response  $\pi_{\max} \in \Pi_{\max}$  satisfying:
        $\mu(v) = \mathbf{r}(v, \pi_{\min}(v)) + \gamma\mu(\pi_{\min}(v))$  for all  $v \in V_{\min}$ ,
        $\mu(v) = \max_{(v, v') \in E} (\mathbf{r}(v, v') + \gamma\mu(v'))$  for all  $v \in V_{\max}$ .
7:   // Policy Improvement Phase
8:   Let  $I = \{v \in V_{\min} \mid \exists (v, v') \in E \text{ such that } \mathbf{r}(v, v') + \gamma\mu(v') < \mu(v)\}$ .
9:   if  $I = \emptyset$  then
10:     return  $(\pi_{\min}, \pi_{\max}, \mu)$ 
11:   end if
12:   Update  $\pi_{\min}$  by switching the edge for all vertices  $v \in I$  to minimize the value:
        $\pi_{\min}(v) \leftarrow \arg \min_{(v, v') \in E} (\mathbf{r}(v, v') + \gamma\mu(v'))$ .
13: end loop

```

Algorithm 2 Restarting Policy Iteration: $\text{RePI}(G, \gamma, \mathbf{r})$

```

1: Input: Graph  $G = (V_{\min}, V_{\max}, E)$  with  $n = |V| \geq 2$ , target discount factor  $\gamma \in [0, 1)$ , reward
   vector  $\mathbf{r} \in \mathbb{R}^E$ .
2: Output: Optimal strategy pair  $(\pi_{\min}^*, \pi_{\max}^*)$  and value vector  $\mu^{\gamma, \mathbf{r}}$  for the game  $G^{\gamma, \mathbf{r}}$ .
3: Let  $k_{\max} = \lceil n \ln n \rceil$ .
4: for  $k = 0, 1, \dots, k_{\max}$  do
5:   Set the threshold discount factor  $\gamma_k = 1 - e^{-k}$ .
6:   if  $\gamma \leq \gamma_k$  then
7:     return  $\text{PI}(G, \gamma, \mathbf{r})$ 
8:   else
9:     // Phase 1: Solve at the threshold
10:     $(\pi_{\min}, \pi_{\max}, \mu_{\gamma_k}) \leftarrow \text{PI}(G, \gamma_k, \mathbf{r})$ 
11:    // Phase 2: Optimality test for the target discount factor
12:    Compute the value vector  $\mu_{\gamma} \in \mathbb{R}^V$  by evaluating  $(\pi_{\min}, \pi_{\max})$  under  $\gamma$ .
13:    Let  $I_{\min} = \{v \in V_{\min} \mid \exists (v, v') \in E \text{ such that } \mathbf{r}(v, v') + \gamma\mu_{\gamma}(v') < \mu_{\gamma}(v)\}$ .
14:    Let  $I_{\max} = \{v \in V_{\max} \mid \exists (v, v') \in E \text{ such that } \mathbf{r}(v, v') + \gamma\mu_{\gamma}(v') > \mu_{\gamma}(v)\}$ .
15:    if  $I_{\min} = \emptyset$  and  $I_{\max} = \emptyset$  then
16:      return  $(\pi_{\min}, \pi_{\max}, \mu_{\gamma})$ 
17:    end if
18:  end if
19: end for
20: // Phase 3: Fallback
21: return  $\text{PI}(G, \gamma, \mathbf{r})$ 

```

537 C Proofs of Section 6

538 In this section, we provide the proofs omitted from Section 6, organized by the same ingredient
539 structure as the main body.

540 C.1 Ingredient 1: Value Decomposition and Lipschitz Continuity

541 **Lemma 2** (Value decomposition under Blackwell optimality). *Consider a game graph G with*
542 *reward vector $\mathbf{r} \in \mathbb{R}^E$. Let $v \in V$ be a vertex.*

543 1. Decomposition: *For all $\gamma \in (\gamma_{bw}(\mathbf{r}), 1)$, the discounted-sum value admits the following decom-*
544 *position, where $w_v^{\mathbf{r}}(\gamma)$ is a rational function of γ :*

$$\mu^{\gamma, \mathbf{r}}(v) = \frac{\mu^{mp, \mathbf{r}}(v)}{1 - \gamma} + w_v^{\mathbf{r}}(\gamma).$$

545 2. Magnitude and derivative: *For all $\gamma \in (\gamma_{bw}(\mathbf{r}), 1)$, the magnitude and the derivative of $w_v^{\mathbf{r}}(\gamma)$*
546 *are bounded by:*

$$|w_v^{\mathbf{r}}(\gamma)| \leq 4n \|\mathbf{r}\|_{\infty} \quad \text{and} \quad \left| \left(\frac{d}{d\gamma} w_v^{\mathbf{r}}(\gamma) \right) \right| \leq 7n^2 \|\mathbf{r}\|_{\infty}.$$

547 *Proof.* Let $M = \|\mathbf{r}\|_{\infty}$ and $\mu \triangleq \mu^{mp, \mathbf{r}}(v)$. Since $\gamma \in (\gamma_{bw}(\mathbf{r}), 1)$, we consider a Blackwell optimal
548 pair of positional strategies $(\pi_{\min}^{bw}, \pi_{\max}^{bw})$ for $G^{\gamma, \mathbf{r}}$. Observe that the path from vertex v under the
549 strategy profile $(\pi_{\min}^{bw}, \pi_{\max}^{bw})$ consists of a finite path of length $t \leq n$ and a recurrent cycle of length
550 $c \leq n$. Let r_k denote the reward at step k . Note that the mean-payoff is the cycle's average reward,
551 i.e., $\mu = \frac{1}{c} \sum_{j=0}^{c-1} r_{t+j}$. Therefore, we have $|\mu| \leq M$ and $\sum_{j=0}^{c-1} (r_{t+j} - \mu) = 0$. We now present the
552 two items of the proof.

553 *Value decomposition.* The value vector expands as follows:

$$\begin{aligned} \mu^{\gamma, \mathbf{r}}(v) &= \sum_{j=0}^{t-1} r_j \gamma^j + \frac{\gamma^t}{1 - \gamma^c} \sum_{j=0}^{c-1} r_{t+j} \gamma^j \\ &= \sum_{j=0}^{t-1} (r_j - \mu) \gamma^j + \frac{\gamma^t}{1 - \gamma^c} \sum_{j=0}^{c-1} (r_{t+j} - \mu) \gamma^j + \mu \sum_{j=0}^{t-1} \gamma^j + \frac{\gamma^t \mu}{1 - \gamma^c} \sum_{j=0}^{c-1} \gamma^j \\ &= \sum_{j=0}^{t-1} (r_j - \mu) \gamma^j + \frac{\gamma^t}{1 - \gamma^c} \sum_{j=0}^{c-1} (r_{t+j} - \mu) (\gamma^j - 1) + \mu \left(\frac{1 - \gamma^t}{1 - \gamma} + \frac{\gamma^t}{1 - \gamma} \right) \\ &= \underbrace{\sum_{j=0}^{t-1} (r_j - \mu) \gamma^j + \sum_{j=1}^{c-1} \left(\frac{\sum_{i=0}^{j-1} \gamma^i}{\sum_{i=0}^{c-1} \gamma^i} \right) (\mu - r_{t+j}) \gamma^t}_{\triangleq w_v^{\mathbf{r}}(\gamma)} + \frac{\mu}{1 - \gamma}, \end{aligned}$$

554 where the first equality is due to decomposing the path into a transient path and a cycle, the second
555 equality rearranges the terms to separate the mean-payoff component, the third equality introduces
556 -1 into the cycle sum because $\sum_{j=0}^{c-1} (r_{t+j} - \mu) = 0$, and the final equality factors out $1 - \gamma$. We
557 define $f_j(\gamma) \triangleq \sum_{i=0}^{j-1} \gamma^i$. Then, we rewrite $w_v^{\mathbf{r}}(\gamma)$ as:

$$\begin{aligned} w_v^{\mathbf{r}}(\gamma) &= \sum_{j=0}^{t-1} (r_j - \mu) \gamma^j + \frac{1}{f_c(\gamma)} \sum_{j=1}^{c-1} (\mu - r_{t+j}) \gamma^t f_j(\gamma) \\ &= \frac{\sum_{j=0}^{t-1} (r_j - \mu) \gamma^j f_c(\gamma) + \sum_{j=1}^{c-1} (\mu - r_{t+j}) \gamma^t f_j(\gamma)}{f_c(\gamma)}. \end{aligned}$$

558 Therefore, $w_v^{\mathbf{r}}(\gamma)$ is a rational function of γ , which proves the first item of the statement. We now
 559 bound its magnitude and derivative.

560 *Magnitude bound.* For magnitude, we get

$$\begin{aligned}
 |w_v^{\mathbf{r}}(\gamma)| &= \left| \sum_{j=0}^{t-1} (r_j - \mu) \gamma^j + \sum_{j=1}^{c-1} (\mu - r_{t+j}) \gamma^t \frac{f_j(\gamma)}{f_c(\gamma)} \right| \\
 &\leq \sum_{j=0}^{t-1} |r_j - \mu| \gamma^j + \sum_{j=1}^{c-1} |\mu - r_{t+j}| \gamma^t \frac{f_j(\gamma)}{f_c(\gamma)} \\
 &\leq 2M \sum_{j=0}^{t-1} \gamma^j + 2M \sum_{j=1}^{c-1} \gamma^t \frac{f_j(\gamma)}{f_c(\gamma)} \\
 &\leq 2M \sum_{j=0}^{t-1} 1 + 2M \sum_{j=1}^{c-1} \frac{f_j(\gamma)}{f_c(\gamma)} \\
 &\leq 4Mn = 4n \|\mathbf{r}\|_{\infty},
 \end{aligned}$$

561 where the first inequality is due to the triangle inequality, the second inequality is due to the bounds
 562 $|r_j - \mu| \leq 2M$ and $|\mu - r_{t+j}| \leq 2M$, the third inequality is due to $\gamma^j, \gamma^t \leq 1$ for all $\gamma \in [0, 1]$, the
 563 fourth inequality is due to $f_j(\gamma) \leq f_c(\gamma)$ for all $j < c$, and the final inequality is due to $t, c \leq n$.

564 *Derivative bound.* For the derivative, we first compute $\left| \frac{d}{d\gamma} \frac{f_j(\gamma)}{f_c(\gamma)} \right|$:

$$\begin{aligned}
 \left| \frac{d}{d\gamma} \frac{f_j(\gamma)}{f_c(\gamma)} \right| &= \left| \frac{(\frac{d}{d\gamma} f_j)(\gamma) f_c(\gamma) - f_j(\gamma) (\frac{d}{d\gamma} f_c)(\gamma)}{f_c(\gamma)^2} \right| \\
 &= \left| \frac{(\frac{d}{d\gamma} f_j)(\gamma)}{f_c(\gamma)} - \frac{f_j(\gamma)}{f_c(\gamma)} \cdot \frac{(\frac{d}{d\gamma} f_c)(\gamma)}{f_c(\gamma)} \right| \\
 &\leq \left| \frac{(\frac{d}{d\gamma} f_j)(\gamma)}{f_c(\gamma)} \right| + \left| \frac{(\frac{d}{d\gamma} f_c)(\gamma)}{f_c(\gamma)} \right| \\
 &\leq 2n,
 \end{aligned}$$

565 where the first equality is due to the quotient rule, the second equality is due to rearranging the
 566 terms, the first inequality is due to the triangle inequality and the fact that $|\frac{f_j(\gamma)}{f_c(\gamma)}| \leq 1$, and the
 567 second inequality is due to $(\frac{d}{d\gamma} f_j)(\gamma) = \sum_{i=1}^{j-1} i \gamma^{i-1} \leq \sum_{i=1}^{c-1} i \gamma^{i-1} = (\frac{d}{d\gamma} f_c)(\gamma) \leq c \sum_{i=1}^{c-1} \gamma^{i-1} \leq$
 568 $n f_c(\gamma)$. Now, we bound the derivative of $w_v^{\mathbf{r}}(\gamma)$ as follows:

$$\begin{aligned}
 \left| \frac{d}{d\gamma} w_v^{\mathbf{r}}(\gamma) \right| &\leq \sum_{j=0}^{t-1} |r_j - \mu| \left| \frac{d}{d\gamma} \gamma^j \right| + \sum_{j=1}^{c-1} |\mu - r_{t+j}| \left| \frac{d}{d\gamma} \left(\gamma^t \frac{f_j(\gamma)}{f_c(\gamma)} \right) \right| \\
 &\leq 2M \sum_{j=0}^{t-1} \left| \frac{d}{d\gamma} \gamma^j \right| + 2M \sum_{j=1}^{c-1} \left| \frac{d}{d\gamma} \left(\gamma^t \frac{f_j(\gamma)}{f_c(\gamma)} \right) \right| \\
 &\leq 2M \sum_{j=1}^{t-1} j + 2M \sum_{j=1}^{c-1} \left| \frac{d}{d\gamma} \left(\gamma^t \frac{f_j(\gamma)}{f_c(\gamma)} \right) \right| \\
 &\leq Mn^2 + 2M \sum_{j=1}^{c-1} 3n \\
 &\leq 7n^2 M = 7n^2 \|\mathbf{r}\|_{\infty},
 \end{aligned}$$

569 where the first inequality is due to the triangle inequality, the second inequality is due to $|r_j - \mu| \leq$
570 $2M$ and $|\mu - r_{t+j}| \leq 2M$, the third inequality is due to $\left| \frac{d}{d\gamma} \gamma^j \right| = j\gamma^{j-1} \leq j$ for all $\gamma \in [0, 1]$, the
571 fourth inequality is due to $t \leq n$ and $\left| \frac{d}{d\gamma} \frac{f_j(\gamma)}{f_c(\gamma)} \right| \leq 2n$ as shown above, and the final inequality is due
572 to $c \leq n$.

573 Hence all the desired items are proved and the result follows. \square

574 **Lemma 3** (Lipschitz continuity). *Let $\mathbf{r} \in \mathbb{R}^E$ be any reward vector. For any edge $a = (v, v') \in E$*
575 *where $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$, the map $\gamma \mapsto \Delta^{\gamma, \mathbf{r}}(v, v')$ is $18n^2 \|\mathbf{r}\|_\infty$ -Lipschitz continuous over the*
576 *interval $(\gamma_{\text{bw}}(\mathbf{r}), 1)$.*

577 *Proof.* By Lemma 2, for all $\gamma \in [\gamma_{\text{bw}}(\mathbf{r}), 1)$, the value at vertices $v, v' \in V$ decomposes as:

$$\mu^{\gamma, \mathbf{r}}(v) = \frac{\mu^{\text{mp}, \mathbf{r}}(v)}{1 - \gamma} + w_v^{\mathbf{r}}(\gamma), \quad \text{and} \quad \mu^{\gamma, \mathbf{r}}(v') = \frac{\mu^{\text{mp}, \mathbf{r}}(v')}{1 - \gamma} + w_{v'}^{\mathbf{r}}(\gamma),$$

578 where $w_v^{\mathbf{r}}(\gamma)$ and $w_{v'}^{\mathbf{r}}(\gamma)$ are rational functions of γ with bounded magnitude and derivative as
579 stated in Lemma 2.

580 Substituting this decomposition into the suboptimality gap yields:

$$\begin{aligned} \Delta^{\gamma, \mathbf{r}}(v, v') &= \mathbf{r}(v, v') + \gamma \left(\frac{\mu^{\text{mp}, \mathbf{r}}(v')}{1 - \gamma} + w_{v'}^{\mathbf{r}}(\gamma) \right) - \left(\frac{\mu^{\text{mp}, \mathbf{r}}(v)}{1 - \gamma} + w_v^{\mathbf{r}}(\gamma) \right) \\ &= \mathbf{r}(v, v') - \mu^{\text{mp}, \mathbf{r}}(v) + \gamma w_{v'}^{\mathbf{r}}(\gamma) - w_v^{\mathbf{r}}(\gamma), \end{aligned}$$

581 where the first equality is due to substituting the value decomposition, and the second equality is
582 due to the assumption that $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$. Notice that the singular term $\frac{\mu^{\text{mp}, \mathbf{r}}(v)}{1 - \gamma}$ exactly
583 cancels out. To bound the Lipschitz constant over $[\gamma_{\text{bw}}(\mathbf{r}), 1)$, we evaluate the derivative with
584 respect to γ :

$$\frac{d}{d\gamma} \Delta^{\gamma, \mathbf{r}}(v, v') = w_{v'}^{\mathbf{r}}(\gamma) + \gamma \left(\frac{d}{d\gamma} w_{v'}^{\mathbf{r}}(\gamma) \right) - \left(\frac{d}{d\gamma} w_v^{\mathbf{r}}(\gamma) \right).$$

585 Therefore, we obtain:

$$\begin{aligned} \left| \frac{d}{d\gamma} \Delta^{\gamma, \mathbf{r}}(v, v') \right| &\leq |w_{v'}^{\mathbf{r}}(\gamma)| + \left| \left(\frac{d}{d\gamma} w_{v'}^{\mathbf{r}}(\gamma) \right) \right| + \left| \left(\frac{d}{d\gamma} w_v^{\mathbf{r}}(\gamma) \right) \right| \\ &\leq 4n \|\mathbf{r}\|_\infty + 7n^2 \|\mathbf{r}\|_\infty + 7n^2 \|\mathbf{r}\|_\infty \leq 18n^2 \|\mathbf{r}\|_\infty, \end{aligned}$$

586 where the first inequality is due to the triangle inequality and the fact that $\gamma < 1$, and the second
587 inequality is due to the bounds on the magnitude and derivative of $w_v^{\mathbf{r}}(\gamma)$ and $w_{v'}^{\mathbf{r}}(\gamma)$ as stated in
588 Lemma 2. Hence, the map $\gamma \mapsto \Delta^{\gamma, \mathbf{r}}(v, v')$ is Lipschitz continuous over $[\gamma_{\text{bw}}(\mathbf{r}), 1)$ with Lipschitz
589 constant $18n^2 \|\mathbf{r}\|_\infty$. \square

590 We conclude by showing that under the smoothed model, with probability one, every vertex has
591 a unique outgoing edge with zero suboptimality gap (Proposition 12), which yields the uniqueness of
592 the Blackwell optimal strategy pair (Corollary 13). These results use standard ideas from smoothed
593 analysis. The following arguments exploit the bounded density of the noise by conditioning on all
594 but one coordinate at a time, which calls for notation that isolates a single edge.

595 **Single-edge decomposition.** To analyze the effect of perturbing a single edge, we decompose
596 the reward and noise vectors along any edge $a \in E$. We write

$$\mathbf{r} = (\mathbf{r}(a); \mathbf{r}_{-a}) \quad \text{and} \quad \xi = (\xi_a; \xi_{-a}),$$

597 where $\mathbf{r}_{-a} \in \mathbb{R}^{E \setminus \{a\}}$ is the vector of rewards for all edges except a , defined by $\mathbf{r}_{-a}(v, v') = \mathbf{r}(v, v')$
598 for all $(v, v') \in E \setminus \{a\}$, and ξ_{-a} is defined analogously. For any fixed $\mathbf{r}_{-a} \in \mathbb{R}^{E \setminus \{a\}}$, we write $\mathbb{P}_{\xi_a}[\cdot]$
599 for the probability measure over the single noise coordinate ξ_a .

600 **Proposition 12** (Optimality uniqueness). *Under the smoothed model, we have*

$$\mathbb{P}_\xi[\forall v \in V, \exists! a = (v, v') \in E \text{ such that } \mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v') \text{ and } \Delta^{1, \mathbf{r}}(a) = 0] = 1.$$

601 *In words, with probability 1, for each vertex $v \in V$, there exists a unique outgoing edge $a = (v, v')$*
 602 *such that $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$ and $\Delta^{1, \mathbf{r}}(a) = 0$.*

603 *Proof.* We separately establish (i) existence, i.e., for every vertex $v \in V$ at least one such edge
 604 exists deterministically; and (ii) uniqueness, i.e., with probability 1 no two distinct outgoing edges
 605 from v both satisfy the conditions. A union bound over the finite vertex set V completes the proof.

606 *Existence.* Let $\mathbf{r} \in \mathbb{R}^E$ be any reward vector and $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ be a Blackwell optimal strategy
 607 pair. Fix a vertex $v \in V$ and let $a = (v, v') \in E$ be the edge selected at v by this strategy pair.
 608 Note that the path from v under the strategy profile $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ consists of a finite path of length
 609 $t \leq n$ and a recurrent cycle of length $c \leq n$. The analogous path from v' contains the same cycle.
 610 Since the mean-payoff is the average reward of the cycle and Blackwell optimal strategies are also
 611 mean-payoff optimal [BK76], we have $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$. Moreover, since a is an optimal action,
 612 we have $\Delta^{\gamma, \mathbf{r}}(a) = 0$ for all $\gamma \in (\gamma_{\text{bw}}(\mathbf{r}), 1)$. By Lemma 3, the limit extends to $\Delta^{1, \mathbf{r}}(a) = 0$, which
 613 yields the existence.

614 *Uniqueness.* Fix a vertex $v \in V$ with at least two outgoing edges, otherwise uniqueness holds. Let
 615 $a_1 = (v, v_1)$ and $a_2 = (v, v_2)$ be any two distinct outgoing edges from v . We define the event \mathcal{E} that
 616 both a_1 and a_2 satisfy the conditions:

$$\mathcal{E} \triangleq \{\mathbf{r} \in \mathbb{R}^E : \mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v_1) = \mu^{\text{mp}, \mathbf{r}}(v_2), \Delta^{1, \mathbf{r}}(a_1) = 0, \Delta^{1, \mathbf{r}}(a_2) = 0\},$$

617 where the mean-payoff equality at v_2 is included so that $\Delta^{1, \mathbf{r}}(a_2)$ is well-defined. We show that
 618 $\mathbb{P}_\xi[\mathcal{E}] = 0$, which implies that with probability 1 at most one outgoing edge from v satisfies the
 619 conditions. Since $a_1 \neq a_2$, a Blackwell optimal strategy must differ from at least one of a_1 or a_2 at
 620 vertex v . Therefore, we decompose $\mathcal{E} \subseteq \mathcal{E}_1 \cup \mathcal{E}_2$, where

$$\begin{aligned} \mathcal{E}_1 &\triangleq \mathcal{E} \cap \{\mathbf{r} \in \mathbb{R}^E : \exists (\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}}) \text{ Blackwell optimal for } \mathbf{r} \text{ with } \pi^{\text{bw}}(v) \neq v_1\}, \\ \mathcal{E}_2 &\triangleq \mathcal{E} \cap \{\mathbf{r} \in \mathbb{R}^E : \exists (\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}}) \text{ Blackwell optimal for } \mathbf{r} \text{ with } \pi^{\text{bw}}(v) \neq v_2\}, \end{aligned}$$

621 and $\pi^{\text{bw}}(v) \triangleq \pi_{\min}^{\text{bw}}(v)$ if $v \in V_{\min}$ and $\pi^{\text{bw}}(v) \triangleq \pi_{\max}^{\text{bw}}(v)$ if $v \in V_{\max}$. We show that $\mathbb{P}_\xi[\mathcal{E}_1] = 0$.
 622 The bound $\mathbb{P}_\xi[\mathcal{E}_2] = 0$ follows by the symmetric argument with the roles of a_1 and a_2 swapped.
 623 Consider any fixed $\mathbf{r}_{-a_1} \in \mathbb{R}^{E \setminus \{a_1\}}$. We analyze the event \mathcal{E}_1 with respect to the random variable
 624 $x \triangleq \mathbf{r}(a_1)$. Formally, we prove $\mathbb{P}_{\xi_{a_1}}[\mathcal{E}_1] = 0$. Let $\mu^{\gamma, \mathbf{r}^{-a_1}}$ denote the optimal value vector of
 625 the restricted game played on the edge set $E \setminus \{a_1\}$. Since this restricted game is completely
 626 independent of a_1 , its value vector $\mu^{\gamma, \mathbf{r}^{-a_1}}$ is a function depending only on \mathbf{r}_{-a_1} . Indeed, for any
 627 $x \in \mathbb{R}$ such that $\mathbf{r} \triangleq (x, \mathbf{r}_{-a_1}) \in \mathcal{E}_1$, there exists a Blackwell optimal strategy pair $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ where
 628 $\pi^{\text{bw}}(v) \neq v_1$. Since the edge $(v, \pi^{\text{bw}}(v)) \in E \setminus \{a_1\}$, removing a_1 does not change the optimal value
 629 at vertex v . Consequently, the value vectors of the full game and the restricted game are equal,
 630 i.e., $\mu^{\gamma, \mathbf{r}} = \mu^{\gamma, \mathbf{r}^{-a_1}}$ for all $\gamma \in (\gamma_{\text{bw}}(\mathbf{r}), 1)$. This equality also implies $\mu^{\text{mp}, \mathbf{r}^{-a_1}}(v) = \mu^{\text{mp}, \mathbf{r}^{-a_1}}(v_1)$ by
 631 Eq. (2). Under the event \mathcal{E}_1 , taking the limit of suboptimality gap evaluates to:

$$\Delta^{1, \mathbf{r}}(a_1) = x + \lim_{\gamma \uparrow 1} (\gamma \mu^{\gamma, \mathbf{r}^{-a_1}}(v_1) - \mu^{\gamma, \mathbf{r}^{-a_1}}(v)).$$

632 Let $C(\mathbf{r}_{-a_1})$ denote the limit term. By Lemma 3 applied to the restricted game, $C(\mathbf{r}_{-a_1})$ exists and is
 633 finite. Since $C(\mathbf{r}_{-a_1})$ is fixed with respect to x , the condition $\Delta^{1, \mathbf{r}}(a_1) = 0$ requires $x + C(\mathbf{r}_{-a_1}) = 0$,
 634 or equivalently, $x = -C(\mathbf{r}_{-a_1})$. Let $\lambda(\mathcal{E}_1)$ denote the Lebesgue measure of the set of x values
 635 satisfying this condition. Note that $\lambda(\mathcal{E}_1) = 0$. Since the probability density of x is uniformly
 636 bounded by ϕ , we have $\mathbb{P}_{\xi_{a_1}}[\mathcal{E}_1] = 0$. Taking the expectation over the independent randomness in

637 ξ_{-a_1} gives

$$\mathbb{P}_\xi[\mathcal{E}_1] = \mathbb{E}_{\xi_{-a_1}} [\mathbb{P}_{\xi_{a_1}}[\mathcal{E}_1]] = 0.$$

638 The symmetric argument applied to a_2 gives $\mathbb{P}_\xi[\mathcal{E}_2] = 0$, so a union bound yields $\mathbb{P}_\xi[\mathcal{E}] \leq \mathbb{P}_\xi[\mathcal{E}_1] +$
 639 $\mathbb{P}_\xi[\mathcal{E}_2] = 0$. A further union bound over the finitely many pairs of distinct outgoing edges at every
 640 vertex $v \in V$ completes the proof. \square

641 **Corollary 13** (Uniqueness of Blackwell optimal strategies). *Under the smoothed model, we have*

$$\mathbb{P}_\xi[\exists! (\pi_{\min}^{bw}, \pi_{\max}^{bw}) \in \Pi_{\min} \times \Pi_{\max} \text{ that is Blackwell optimal for } \mathbf{r}] = 1.$$

642 *In words, with probability 1, there exists a unique Blackwell optimal strategy pair.*

643 *Proof.* Existence is guaranteed deterministically by the classical Blackwell optimality result [Bla62,
 644 BK76]. It remains to establish uniqueness with probability 1. By Proposition 12, the event

$$\mathcal{E} \triangleq \{\mathbf{r} \in \mathbb{R}^E : \forall v \in V, \exists! a^*(v) = (v, v') \in E \text{ s.t. } \mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v') \text{ and } \Delta^{1, \mathbf{r}}(a^*(v)) = 0\}$$

645 satisfies $\mathbb{P}_\xi[\mathcal{E}] = 1$. We show that on \mathcal{E} , the Blackwell optimal strategy pair is unique.

646 Fix $\mathbf{r} \in \mathcal{E}$ and let $(\pi_{\min}^{bw}, \pi_{\max}^{bw})$ be any Blackwell optimal strategy pair for \mathbf{r} . For any vertex
 647 $v \in V$, let $a = (v, \pi^{bw}(v)) \in E$ be the edge selected at v by this pair, where $\pi^{bw}(v) \triangleq \pi_{\min}^{bw}(v)$ if
 648 $v \in V_{\min}$ and $\pi^{bw}(v) \triangleq \pi_{\max}^{bw}(v)$ if $v \in V_{\max}$. Since Blackwell optimal strategies are also mean-payoff
 649 optimal [BK76], we have $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(\pi^{bw}(v))$. Moreover, since Blackwell optimal strategies
 650 are optimal in discounted-sum games for all $\gamma \in (\gamma_{bw}(\mathbf{r}), 1)$, we have $\Delta^{\gamma, \mathbf{r}}(a) = 0$ for such γ , and
 651 by Lemma 3, taking the limit yields $\Delta^{1, \mathbf{r}}(a) = 0$. Since $\mathbf{r} \in \mathcal{E}$, the unique edge satisfying these
 652 conditions at v is $a^*(v)$, so $a = a^*(v)$. Since this holds for every $v \in V$, both π_{\min}^{bw} and π_{\max}^{bw} are
 653 uniquely determined, which completes the proof. \square

654 **Lemma 14** (Unique optimal strategies at fixed thresholds). *Let $\Gamma \subseteq [0, 1)$ be any finite set of*
 655 *discount factors. Under the smoothed model, with probability 1, for every $\gamma \in \Gamma$ with $\gamma > \gamma_{bw}(\mathbf{r})$, the*
 656 *discounted-sum game $G^{\gamma, \mathbf{r}}$ has a unique optimal strategy pair, and this pair is the unique Blackwell*
 657 *optimal strategy pair for \mathbf{r} .*

658 *Proof.* Fix a discount factor $\gamma \in [0, 1)$. We first show that, with probability 1, the game $G^{\gamma, \mathbf{r}}$ has
 659 a unique optimal strategy pair. It is enough to show that, with probability 1, no vertex has two
 660 distinct outgoing edges that both satisfy the discounted Bellman equality $\Delta^{\gamma, \mathbf{r}}(a) = 0$. Indeed, by
 661 the Bellman optimality equations, at every vertex at least one outgoing edge satisfies this equality,
 662 and if exactly one such edge exists at every vertex, then both players' optimal positional strategies
 663 are uniquely determined.

664 Fix a vertex v with two distinct outgoing edges $a_1 = (v, v_1)$ and $a_2 = (v, v_2)$, and consider the
 665 event

$$\mathcal{E} \triangleq \{\mathbf{r} \in \mathbb{R}^E : \Delta^{\gamma, \mathbf{r}}(a_1) = 0 \text{ and } \Delta^{\gamma, \mathbf{r}}(a_2) = 0\}.$$

666 We prove that $\mathbb{P}_\xi[\mathcal{E}] = 0$. Condition on all rewards except a_1 , writing $x = \mathbf{r}(a_1)$ and $\mathbf{r} = (x; \mathbf{r}_{-a_1})$.
 667 On the event \mathcal{E} , the edge a_2 is also optimal at v ; hence there exists an optimal strategy pair for
 668 $G^{\gamma, \mathbf{r}}$ whose strategy at v selects a_2 and therefore does not use a_1 . Removing a_1 from the game
 669 does not change the value vector at discount γ : if $v \in V_{\min}$, Min still has an optimal strategy that
 670 avoids a_1 , while removing a Min edge can only increase the value; if $v \in V_{\max}$, Max still has an
 671 optimal strategy that avoids a_1 , while removing a Max edge can only decrease the value. Thus
 672 $\mu^{\gamma, \mathbf{r}} = \mu^{\gamma, \mathbf{r}_{-a_1}}$ on \mathcal{E} , where $\mu^{\gamma, \mathbf{r}_{-a_1}}$ denotes the value vector of the game with edge a_1 removed.
 673 Consequently, on \mathcal{E} ,

$$0 = \Delta^{\gamma, \mathbf{r}}(a_1) = x + \gamma \mu^{\gamma, \mathbf{r}_{-a_1}}(v_1) - \mu^{\gamma, \mathbf{r}_{-a_1}}(v).$$

674 For the fixed conditioned vector \mathbf{r}_{-a_1} , the last two terms are independent of x , so this equality
675 can hold for at most one value of x . Since the conditional density of $x = \mathbf{r}(a_1)$ is bounded, the
676 conditional probability of \mathcal{E} is zero. Integrating over \mathbf{r}_{-a_1} gives $\mathbb{P}_\xi[\mathcal{E}] = 0$. A union bound over all
677 vertices and all pairs of distinct outgoing edges shows that $G^{\gamma, \mathbf{r}}$ has a unique optimal strategy pair
678 with probability 1.

679 Since Γ is finite, another union bound implies that, with probability 1, every game $G^{\gamma, \mathbf{r}}$ with
680 $\gamma \in \Gamma$ has a unique optimal strategy pair. Intersect this event with the probability-one event of
681 [Corollary 13](#). On the intersection, if $\gamma \in \Gamma$ and $\gamma > \gamma_{\text{bw}}(\mathbf{r})$, then the unique Blackwell optimal pair
682 is optimal for $G^{\gamma, \mathbf{r}}$ by the definition of the Blackwell threshold. Since the discounted optimal pair
683 at γ is unique, it must coincide with the unique Blackwell optimal pair. \square

684 C.2 Ingredient 2: Suboptimality Gap Separation

685 **Lemma 4** (Suboptimality separation). *Let G be any game graph and let $a = (v, v') \in E$ be any*
686 *fixed edge. Under the smoothed model, for any $\delta > 0$, we have*

$$\mathbb{P}_\xi[\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v') \text{ and } 0 < |\Delta^{1, \mathbf{r}}(a)| < \delta] \leq 2\delta\phi.$$

687 *Proof.* If a is the only outgoing edge from v , it must be selected by any valid strategy. Thus,
688 $\Delta^{1, \mathbf{r}}(a) = 0$ whenever well-defined, and the event has probability 0. Assume a is not the only
689 outgoing edge from v . Consider any fixed $\mathbf{r}_{-a} \in \mathbb{R}^{E \setminus \{a\}}$. We bound the probability of the event

$$\mathcal{E} \triangleq \{\mathbf{r}(a) \in \mathbb{R} : \mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v') \text{ and } 0 < |\Delta^{1, \mathbf{r}}(a)| < \delta\}.$$

690 Formally, we prove $\mathbb{P}_{\xi_a}[\mathcal{E}] \leq 2\delta\phi$.

691 Let $\mu^{\gamma, \mathbf{r}^{-a}}$ denote the optimal value vector of the restricted game played on the edge set $E \setminus \{a\}$.
692 Since this restricted game is completely independent of edge a , its value vector $\mu^{\gamma, \mathbf{r}^{-a}}$ is a function
693 depending only on \mathbf{r}_{-a} . For any $x \in \mathcal{E}$, let $\mathbf{r} = (x, \mathbf{r}_{-a})$. Then, we have $\Delta^{1, \mathbf{r}}(a) \neq 0$. This implies
694 that the edge a is strictly suboptimal under the Blackwell optimal strategies. Consequently, for γ
695 sufficiently close to 1, the optimal strategy does not select a , meaning the value vectors of the full
696 game and the restricted game exactly match: $\mu^{\gamma, \mathbf{r}} = \mu^{\gamma, \mathbf{r}^{-a}}$. Taking the limit $\gamma \uparrow 1$ of $(1 - \gamma)\mu^{\gamma, \mathbf{r}}$
697 and applying [Eq. \(2\)](#) gives $\mu^{\text{mp}, \mathbf{r}^{-a}}(u) = \mu^{\text{mp}, \mathbf{r}}(u)$ for all $u \in V$. In particular, the assumption
698 $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$ implies $\mu^{\text{mp}, \mathbf{r}^{-a}}(v) = \mu^{\text{mp}, \mathbf{r}^{-a}}(v')$.

699 Under this condition, the limit of the suboptimality gap evaluates to:

$$\Delta^{1, \mathbf{r}}(a) = x + \lim_{\gamma \uparrow 1} (\gamma \mu^{\gamma, \mathbf{r}^{-a}}(v') - \mu^{\gamma, \mathbf{r}^{-a}}(v)).$$

700 Let $C(\mathbf{r}_{-a})$ denote the limit term. Applying the value decomposition ([Lemma 2](#)) to the restricted
701 game and using $\mu^{\text{mp}, \mathbf{r}^{-a}}(v) = \mu^{\text{mp}, \mathbf{r}^{-a}}(v')$, the singular $\frac{1}{1-\gamma}$ terms cancel, so $C(\mathbf{r}_{-a})$ exists and is
702 finite. Since $C(\mathbf{r}_{-a})$ is fixed with respect to x , the condition $0 < |x + C(\mathbf{r}_{-a})| < \delta$ requires the value
703 of x to fall strictly within the interval $(-C(\mathbf{r}_{-a}) - \delta, -C(\mathbf{r}_{-a}) + \delta)$. Let $\lambda(\mathcal{E})$ denote the Lebesgue
704 measure of the set of x values satisfying this condition. Note that $\lambda(\mathcal{E}) \leq 2\delta$. Since the probability
705 density of x is uniformly bounded by ϕ , the probability of the event satisfies $\mathbb{P}_{\xi_a}[\mathcal{E}] \leq 2\delta\phi$. Taking
706 expectation over the independent randomness in ξ_{-a} gives

$$\mathbb{P}_\xi[\mathcal{E}] = \mathbb{E}_{\xi_{-a}}[\mathbb{P}_{\xi_a}[\mathcal{E}]] \leq 2\delta\phi,$$

707 which completes the proof. \square

708 C.3 Ingredient 3: Mean-payoff Value Separation

709 **Lemma 5** (Mean-payoff separation). *Let G be any game graph and let $v, v' \in V$. Under the*
 710 *smoothed model, for any $B > 0$, we have*

$$\mathbb{P}_\xi \left[0 < |\mu^{\text{mp}, \mathbf{r}}(v) - \mu^{\text{mp}, \mathbf{r}}(v')| < B \right] \leq 2Bmn\phi.$$

711 To prove the above result, we first define some useful notation. We then establish a reachability
 712 separation lemma, which is a deterministic statement about the structure of Blackwell optimal
 713 strategies. Finally, we use this lemma to prove the mean-payoff value separation result.

714 *Cycles.* Let \mathcal{C} denote the set of simple directed cycles of G . For a cycle $C \in \mathcal{C}$, let $E(C)$ and $V(C)$
 715 be its edge set and vertex set, and write

$$\text{avg}_{\mathbf{r}}(C) \triangleq \frac{1}{|E(C)|} \sum_{a \in E(C)} \mathbf{r}(a).$$

716 *Reachable vertices and edges.* For a strategy $\pi_{\min} \in \Pi_{\min}$ of the Min Player, let $V_{\pi_{\min}}(v)$ be the
 717 set of vertices reachable from v after fixing Min to π_{\min} and leaving all Max edges available. Let
 718 $E_{\pi_{\min}}(v)$ be the set of edges reachable from v in this one-player game, i.e.,

$$E_{\pi_{\min}}(v) \triangleq \{(v', v'') \in E : v' \in V_{\pi_{\min}}(v) \cap V_{\max}\} \cup \{(v', \pi_{\min}(v')) : v' \in V_{\pi_{\min}}(v) \cap V_{\min}\}.$$

719 Symmetrically, for a strategy $\pi_{\max} \in \Pi_{\max}$ of the Max Player, let $V_{\pi_{\max}}(v)$ be the set of vertices
 720 reachable from v after fixing Max to π_{\max} and leaving all Min edges available, and let

$$E_{\pi_{\max}}(v) \triangleq \{(v', v'') \in E : v' \in V_{\pi_{\max}}(v) \cap V_{\min}\} \cup \{(v', \pi_{\max}(v')) : v' \in V_{\pi_{\max}}(v) \cap V_{\max}\}.$$

721 For a strategy $\pi_{\min} \in \Pi_{\min}$, the cycles $C \in \mathcal{C}$ with $E(C) \subseteq E_{\pi_{\min}}(v)$ are exactly the simple directed
 722 cycles reachable from v after fixing π_{\min} , and likewise for $\pi_{\max} \in \Pi_{\max}$.

723 **Lemma 15** (Reachability separation). *Let $\mathbf{r} \in \mathbb{R}^E$, let $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ be a Blackwell optimal strategy*
 724 *pair, and let $v, v' \in V$. The cycle $C_v^{\mathbf{r}}$ is the recurrent cycle reached from v under $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$.*

725 1. *If $\mu^{\text{mp}, \mathbf{r}}(v) > \mu^{\text{mp}, \mathbf{r}}(v')$, then*

$$E(C_v^{\mathbf{r}}) \cap E_{\pi_{\min}^{\text{bw}}}(v') = \emptyset,$$

726 2. *If $\mu^{\text{mp}, \mathbf{r}}(v) < \mu^{\text{mp}, \mathbf{r}}(v')$, then*

$$E(C_v^{\mathbf{r}}) \cap E_{\pi_{\max}^{\text{bw}}}(v') = \emptyset.$$

727 *Proof.* Since the pair $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ is Blackwell optimal, it is also mean-payoff optimal [BK76].
 728 Therefore, π_{\min}^{bw} is an optimal min strategy and π_{\max}^{bw} an optimal max strategy for the game $G^{\text{mp}, \mathbf{r}}$.
 729 Suppose $E(C_v^{\mathbf{r}}) \cap E_{\pi_{\min}^{\text{bw}}}(v') \neq \emptyset$. Then, there exists a vertex $u \in C_v^{\mathbf{r}} \cap V_{\pi_{\min}^{\text{bw}}}(v')$. Since $C_v^{\mathbf{r}}$ follows
 730 π_{\min}^{bw} at every Min vertex, traversing it from u gives $E(C_v^{\mathbf{r}}) \subseteq E_{\pi_{\min}^{\text{bw}}}(v')$, so some Max strategy π_{\max}
 731 makes $\rho(v', \pi_{\min}^{\text{bw}}, \pi_{\max})$ reach $C_v^{\mathbf{r}}$ and cycle it forever, giving $u^{\text{mp}, \mathbf{r}}(\rho(v', \pi_{\min}^{\text{bw}}, \pi_{\max})) = \text{avg}_{\mathbf{r}}(C_v^{\mathbf{r}})$.
 732 As π_{\min}^{bw} is an optimal min strategy, this payoff is at most $\mu^{\text{mp}, \mathbf{r}}(v')$, and therefore

$$\mu^{\text{mp}, \mathbf{r}}(v) = \text{avg}_{\mathbf{r}}(C_v^{\mathbf{r}}) \leq \mu^{\text{mp}, \mathbf{r}}(v'),$$

733 contradicting $\mu^{\text{mp}, \mathbf{r}}(v) > \mu^{\text{mp}, \mathbf{r}}(v')$. Hence $E(C_v^{\mathbf{r}}) \cap E_{\pi_{\min}^{\text{bw}}}(v') = \emptyset$.

734 The proof of the second statement is symmetric. Suppose $E(C_v^{\mathbf{r}}) \cap E_{\pi_{\max}^{\text{bw}}}(v') \neq \emptyset$. Then, there
 735 exists a vertex $u \in C_v^{\mathbf{r}} \cap V_{\pi_{\max}^{\text{bw}}}(v')$. Since $C_v^{\mathbf{r}}$ follows π_{\max}^{bw} at every Max vertex, traversing it from u
 736 gives $E(C_v^{\mathbf{r}}) \subseteq E_{\pi_{\max}^{\text{bw}}}(v')$, so some Min strategy π_{\min} makes $\rho(v', \pi_{\min}, \pi_{\max}^{\text{bw}})$ reach $C_v^{\mathbf{r}}$ and cycle it
 737 forever, giving $u^{\text{mp}, \mathbf{r}}(\rho(v', \pi_{\min}, \pi_{\max}^{\text{bw}})) = \text{avg}_{\mathbf{r}}(C_v^{\mathbf{r}})$. As π_{\max}^{bw} is an optimal max strategy, this payoff
 738 is at least $\mu^{\text{mp}, \mathbf{r}}(v')$, and therefore

$$\mu^{\text{mp}, \mathbf{r}}(v) = \text{avg}_{\mathbf{r}}(C_v^{\mathbf{r}}) \geq \mu^{\text{mp}, \mathbf{r}}(v'),$$

739 contradicting $\mu^{\text{mp}, \mathbf{r}}(v) < \mu^{\text{mp}, \mathbf{r}}(v')$. Hence $E(C_v^{\mathbf{r}}) \cap E_{\pi_{\max}^{\text{bw}}}(v') = \emptyset$. \square

740 *Proof of Lemma 5.* If $v = v'$, the event is empty; assume $v \neq v'$. Let

$$\mathcal{E}_{\text{bw}} \triangleq \{\mathbf{r} \in \mathbb{R}^E : \exists! (\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}}) \in \Pi_{\min} \times \Pi_{\max} \text{ that is Blackwell optimal for } \mathbf{r}\}.$$

741 By Corollary 13, $\mathbb{P}_\xi[\mathcal{E}_{\text{bw}}] = 1$, so it suffices to bound the two bad events

$$\begin{aligned} \mathcal{B}^+ &\triangleq \mathcal{E}_{\text{bw}} \cap \{\mathbf{r} : 0 < \mu^{\text{mp},\mathbf{r}}(v) - \mu^{\text{mp},\mathbf{r}}(v') < B\}, \\ \mathcal{B}^- &\triangleq \mathcal{E}_{\text{bw}} \cap \{\mathbf{r} : 0 < \mu^{\text{mp},\mathbf{r}}(v') - \mu^{\text{mp},\mathbf{r}}(v) < B\}. \end{aligned}$$

742 We bound $\mathbb{P}_\xi[\mathcal{B}^+]$; the bound on $\mathbb{P}_\xi[\mathcal{B}^-]$ is identical after exchanging v and v' .

743 Given an edge $a \in E$ and rewards \mathbf{r}_{-a} , when $(x; \mathbf{r}_{-a}) \in \mathcal{E}_{\text{bw}}$ let π_{\min}^{bw} be its Blackwell optimal
744 Min strategy and $C_v^{(x; \mathbf{r}_{-a})}$ the recurrent cycle reached from v , and define the slice

$$S_a^+(\mathbf{r}_{-a}) \triangleq \{x \in \mathbb{R} : (x; \mathbf{r}_{-a}) \in \mathcal{E}_{\text{bw}}, 0 < \mu^{\text{mp},(x; \mathbf{r}_{-a})}(v) - \mu^{\text{mp},(x; \mathbf{r}_{-a})}(v') < B, a \in E(C_v^{(x; \mathbf{r}_{-a})}), a \notin E_{\pi_{\min}^{\text{bw}}}(v')\}.$$

745 For rewards $\mathbf{r} \in \mathcal{B}^+$, let $C_v^{\mathbf{r}}$ be the recurrent cycle reached from v under its unique Blackwell optimal
746 pair. Therefore, we have $\text{avg}_{\mathbf{r}}(C_v^{\mathbf{r}}) = \mu^{\text{mp},\mathbf{r}}(v)$. By Lemma 15, for all edges $a \in E(C_v^{\mathbf{r}})$, we have
747 $a \notin E_{\pi_{\min}^{\text{bw}}}(v')$. Thus, $\mathbf{r}(a) \in S_a^+(\mathbf{r}_{-a})$. Consequently, we have

$$\mathcal{B}^+ \subseteq \bigcup_{a \in E} \{\mathbf{r} : \mathbf{r}(a) \in S_a^+(\mathbf{r}_{-a})\}.$$

748 Moreover, we show that $\mu^{\text{mp},(x; \mathbf{r}_{-a})}(v')$ is constant on $S_a^+(\mathbf{r}_{-a})$. For a fixed Min strategy π_{\min} , every
749 play from v' is a lasso whose recurrent cycle C satisfies $E(C) \subseteq E_{\pi_{\min}}(v')$, and the Max Player can
750 reach and repeat any such cycle; hence the largest payoff the Max Player secures from v' against
751 π_{\min} is $\max\{\text{avg}_{(x; \mathbf{r}_{-a})}(C) : C \in \mathcal{C}, E(C) \subseteq E_{\pi_{\min}}(v')\}$. Since the mean-payoff value is the min-max
752 of the payoff,

$$\mu^{\text{mp},(x; \mathbf{r}_{-a})}(v') = \min_{\pi_{\min} \in \Pi_{\min}} \max_{\substack{C \in \mathcal{C} \\ E(C) \subseteq E_{\pi_{\min}}(v')}} \text{avg}_{(x; \mathbf{r}_{-a})}(C).$$

753 For $x \in S_a^+(\mathbf{r}_{-a})$, the Blackwell optimal π_{\min}^{bw} is an optimal min strategy, so it attains this minimum,
754 and it satisfies $a \notin E_{\pi_{\min}^{\text{bw}}}(v')$; hence the minimum is attained among the strategies with $a \notin$
755 $E_{\pi_{\min}}(v')$:

$$\mu^{\text{mp},(x; \mathbf{r}_{-a})}(v') = \min_{\substack{\pi_{\min} \in \Pi_{\min} \\ a \notin E_{\pi_{\min}}(v')}} \max_{\substack{C \in \mathcal{C} \\ E(C) \subseteq E_{\pi_{\min}}(v')}} \text{avg}_{(x; \mathbf{r}_{-a})}(C).$$

756 Every cycle on the right-hand side avoids a , so its average does not involve $x = \mathbf{r}(a)$; hence
757 $\mu^{\text{mp},(x; \mathbf{r}_{-a})}(v')$ is a constant α depending only on \mathbf{r}_{-a} for all $x \in S_a^+(\mathbf{r}_{-a})$. Let $f(x) \triangleq \mu^{\text{mp},(x; \mathbf{r}_{-a})}(v)$.
758 The function f is continuous, piecewise affine, nondecreasing, and absolutely continuous because it
759 is the min-max over finitely many positional strategy pairs, and for each such pair the recurrent-
760 cycle mean is an affine function of x with slope 0 or $1/c$ for some $c \leq n$. Let

$$D \triangleq \bigcup_{\substack{C, C' \in \mathcal{C} \\ \text{avg}_{(\cdot; \mathbf{r}_{-a})}(C) \neq \text{avg}_{(\cdot; \mathbf{r}_{-a})}(C')}} \{x \in \mathbb{R} : \text{avg}_{(x; \mathbf{r}_{-a})}(C) = \text{avg}_{(x; \mathbf{r}_{-a})}(C')\},$$

761 which is finite and, since each affine piece of f coincides with $\text{avg}_{(x; \mathbf{r}_{-a})}(C)$ for some $C \in \mathcal{C}$, contains
762 every breakpoint of f . Fix a connected component J of $\mathbb{R} \setminus D$ with $J \cap S_a^+(\mathbf{r}_{-a}) \neq \emptyset$. Because
763 every breakpoint of f belongs to D , the restriction of f to J consists of a single affine piece. Hence
764 there exists a cycle $\widehat{C} \in \mathcal{C}$ such that

$$f(x) = \text{avg}_{(x; \mathbf{r}_{-a})}(\widehat{C}) \quad \text{for every } x \in J.$$

765 Choose $x_0 \in J \cap S_a^+(\mathbf{r}_{-a})$ and let

$$C_0 \triangleq C_v^{(x_0; \mathbf{r}_{-a})}.$$

766 By the definition of $S_a^+(\mathbf{r}_{-a})$, we have $(x_0; \mathbf{r}_{-a}) \in \mathcal{E}_{\text{bw}}$ and $a \in E(C_0)$. The Blackwell optimal
 767 strategy pair at $(x_0; \mathbf{r}_{-a})$ induces the recurrent cycle C_0 from v , and therefore

$$f(x_0) = \text{avg}_{(x_0; \mathbf{r}_{-a})}(C_0) = \text{avg}_{(x_0; \mathbf{r}_{-a})}(\widehat{C}).$$

768 If the two affine functions $\text{avg}_{(\cdot; \mathbf{r}_{-a})}(C_0)$ and $\text{avg}_{(\cdot; \mathbf{r}_{-a})}(\widehat{C})$ were not identical, their equality at x_0
 769 would imply $x_0 \in D$ by the definition of D , contradicting $x_0 \in J \subseteq \mathbb{R} \setminus D$. Consequently, these
 770 affine functions are identical, and hence

$$f(x) = \text{avg}_{(x; \mathbf{r}_{-a})}(C_0) \quad \text{for every } x \in J.$$

771 The recurrent cycle C_0 is simple, contains a , and has at most n edges. Thus, the variable reward
 772 $x = \mathbf{r}(a)$ occurs exactly once in the sum defining its cycle average, and therefore

$$f'(x) = \frac{1}{|E(C_0)|} \geq \frac{1}{n} \quad \text{for all } x \in J.$$

773 Every $x \in S_a^+(\mathbf{r}_{-a}) \setminus D$ lies in such a component, so $f'(x) \geq 1/n$ for almost every $x \in S_a^+(\mathbf{r}_{-a})$.
 774 Since $\mu^{\text{mp}, (x; \mathbf{r}_{-a})}(v') = \alpha$ on $S_a^+(\mathbf{r}_{-a})$, we have $S_a^+(\mathbf{r}_{-a}) \subseteq f^{-1}((\alpha, \alpha + B))$. Therefore, since f is
 775 nondecreasing and absolutely continuous, we have

$$\lambda(S_a^+(\mathbf{r}_{-a})) \leq n \int_{S_a^+(\mathbf{r}_{-a})} f'(x) dx \leq n \int_{f^{-1}((\alpha, \alpha + B))} f'(x) dx \leq nB,$$

776 where $\lambda(\cdot)$ denotes the Lebesgue measure.

777 Conditioning on \mathbf{r}_{-a} , the reward $\mathbf{r}(a) = \mathbf{r}_0(a) + \xi_a$ has density at most ϕ and is independent of
 778 \mathbf{r}_{-a} , so

$$\mathbb{P}_\xi[\mathcal{B}^+] \leq \sum_{a \in E} \mathbb{P}_\xi[\mathbf{r}(a) \in S_a^+(\mathbf{r}_{-a})] \leq \sum_{a \in E} \phi \mathbb{E}_{\mathbf{r}_{-a}}[\lambda(S_a^+(\mathbf{r}_{-a}))] \leq mnB\phi.$$

779 The identical argument gives $\mathbb{P}_\xi[\mathcal{B}^-] \leq mnB\phi$, and therefore

$$\mathbb{P}_\xi\left[0 < |\mu^{\text{mp}, \mathbf{r}}(v) - \mu^{\text{mp}, \mathbf{r}}(v')| < B\right] \leq \mathbb{P}_\xi[\mathcal{B}^+] + \mathbb{P}_\xi[\mathcal{B}^-] \leq 2mnB\phi.$$

780 □

781 D Proofs of Section 7

782 In this section, we provide the proofs omitted from Section 7, organized by the same structure as
 783 the main body.

784 D.1 Tail Bound

785 **Lemma 6** (Tail bound). *Let G be any game graph. Under the smoothed model, for any $x \geq 1$,*
 786 *with the bounding function defined as $M(x) \triangleq 1 + \frac{1}{\theta} \ln(mx)$, we have*

$$\mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] \leq \frac{KM(x)}{x}.$$

787 *Proof.* For any $M \geq 1$, $B > 0$, and $\delta > 0$, define the events

$$\mathcal{E}_1(M) \triangleq \{\mathbf{r} \in \mathbb{R}^E : \|\mathbf{r}\|_\infty \leq M\},$$

$$\mathcal{E}_2(B) \triangleq \{\mathbf{r} \in \mathbb{R}^E : |\mu^{\text{mp}, \mathbf{r}}(v) - \mu^{\text{mp}, \mathbf{r}}(v')| \in \{0\} \cup [B, \infty) \text{ for all } (v, v') \in E\},$$

$$\mathcal{E}_3(\delta) \triangleq \{\mathbf{r} \in \mathbb{R}^E : |\Delta^{1, \mathbf{r}}(a)| \in \{0\} \cup [\delta, \infty) \text{ for all } a = (v, v') \in E \text{ with } \mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')\},$$

$$\mathcal{E}_4 \triangleq \{\mathbf{r} \in \mathbb{R}^E : \forall v \in V, \exists! a = (v, v') \in E \text{ such that } \mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v') \text{ and } \Delta^{1, \mathbf{r}}(a) = 0\}.$$

788 Let $L \triangleq 18n^2M$. Conditioned on $\mathcal{E}_1(M) \cap \mathcal{E}_2(B) \cap \mathcal{E}_3(\delta) \cap \mathcal{E}_4$, we show that $\gamma_{\text{bw}}(\mathbf{r}) \leq \max(0, 1 -$
789 $\frac{\delta}{2L}, 1 - \frac{B}{10nM}) = \max(0, 1 - \frac{\delta}{36n^2M}, 1 - \frac{B}{10nM})$.

790 Indeed, consider $\mathbf{r} \in \mathcal{E}_1(M) \cap \mathcal{E}_2(B) \cap \mathcal{E}_3(\delta) \cap \mathcal{E}_4$. Since $\mathbf{r} \in \mathcal{E}_4$, by [Proposition 12](#), for each vertex
791 $v \in V$ there exists a unique edge $a^*(v) = (v, v') \in E$ with $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$ and $\Delta^{1, \mathbf{r}}(a^*(v)) = 0$.
792 Therefore, by [Corollary 13](#), there exists a unique Blackwell optimal strategy pair $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$, which
793 selects $a^*(v)$ at each vertex $v \in V$.

794 By [Lemma 3](#) and [Lemma 2](#), since $\mathbf{r} \in \mathcal{E}_1(M)$, the L -Lipschitz bound on the suboptimality
795 gap and the value decomposition both hold on $(\gamma_{\text{bw}}(\mathbf{r}), 1)$. We argue by contradiction: suppose
796 $\gamma_{\text{bw}}(\mathbf{r}) > \max(1 - \frac{\delta}{2L}, 1 - \frac{B}{10nM})$. Then for any $\gamma \in (\gamma_{\text{bw}}(\mathbf{r}), 1)$, $1 - \gamma < 1 - \gamma_{\text{bw}}(\mathbf{r}) < \min(\frac{\delta}{2L}, \frac{B}{10nM})$.
797 We bound $|\Delta^{\gamma, \mathbf{r}}(a)|$ for any non-optimal edge $a = (v, v') \neq a^*(v)$ on $(\gamma_{\text{bw}}(\mathbf{r}), 1)$, splitting into two
798 cases.

799 **Case $\mu^{\text{mp}, \mathbf{r}}(v) = \mu^{\text{mp}, \mathbf{r}}(v')$.** Since $\mathbf{r} \in \mathcal{E}_3(\delta)$ and a is non-optimal, $|\Delta^{1, \mathbf{r}}(a)| \geq \delta$. Then:

$$\begin{aligned} |\Delta^{\gamma, \mathbf{r}}(a)| &\geq |\Delta^{1, \mathbf{r}}(a)| - |\Delta^{\gamma, \mathbf{r}}(a) - \Delta^{1, \mathbf{r}}(a)| \\ &\geq \delta - L(1 - \gamma) > \frac{\delta}{2} > 0, \end{aligned}$$

800 where the first inequality is the triangle inequality, the second combines $|\Delta^{1, \mathbf{r}}(a)| \geq \delta$ with the
801 L -Lipschitz bound on $(\gamma_{\text{bw}}(\mathbf{r}), 1)$, and the strict third inequality uses $1 - \gamma < \frac{\delta}{2L}$.

802 **Case $\mu^{\text{mp}, \mathbf{r}}(v) \neq \mu^{\text{mp}, \mathbf{r}}(v')$.** Since $\mathbf{r} \in \mathcal{E}_2(B)$, $|\mu^{\text{mp}, \mathbf{r}}(v') - \mu^{\text{mp}, \mathbf{r}}(v)| \geq B$. Substituting the value
803 decomposition from [Lemma 2](#) into the suboptimality gap yields

$$\Delta^{\gamma, \mathbf{r}}(a) = \mathbf{r}(v, v') - \mu^{\text{mp}, \mathbf{r}}(v') + \frac{\mu^{\text{mp}, \mathbf{r}}(v') - \mu^{\text{mp}, \mathbf{r}}(v)}{1 - \gamma} + \gamma w_{v'}^{\mathbf{r}}(\gamma) - w_v^{\mathbf{r}}(\gamma).$$

804 Applying the reverse triangle inequality to isolate the singular term:

$$\begin{aligned} |\Delta^{\gamma, \mathbf{r}}(a)| &\geq \frac{|\mu^{\text{mp}, \mathbf{r}}(v') - \mu^{\text{mp}, \mathbf{r}}(v)|}{1 - \gamma} - |\mathbf{r}(v, v') - \mu^{\text{mp}, \mathbf{r}}(v') + \gamma w_{v'}^{\mathbf{r}}(\gamma) - w_v^{\mathbf{r}}(\gamma)| \\ &\geq \frac{B}{1 - \gamma} - (2M + 4nM + 4nM) \\ &\geq \frac{B}{1 - \gamma} - 10nM > 0, \end{aligned}$$

805 where the first inequality is the reverse triangle inequality together with $|\mu^{\text{mp}, \mathbf{r}}(v') - \mu^{\text{mp}, \mathbf{r}}(v)| \geq B$,
806 the second uses $|\mathbf{r}(v, v')| \leq M$ and $|\mu^{\text{mp}, \mathbf{r}}(v')| \leq M$ (from $\mathcal{E}_1(M)$) together with $|w_{v'}^{\mathbf{r}}(\gamma)|, |w_v^{\mathbf{r}}(\gamma)| \leq$
807 $4nM$ from [Lemma 2](#), the third uses $n \geq 1$, and the final strict inequality uses $1 - \gamma < \frac{B}{10nM}$.

808 Combining the two cases, $|\Delta^{\gamma, \mathbf{r}}(a)| > 0$ for every non-optimal edge a and every $\gamma \in (\gamma_{\text{bw}}(\mathbf{r}), 1)$.
809 By definition of Blackwell optimality, on $(\gamma_{\text{bw}}(\mathbf{r}), 1)$ the pair $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ is optimal, so for non-
810 optimal edges a outgoing from v , the sign of $\Delta^{\gamma, \mathbf{r}}(a)$ is positive if $v \in V_{\min}$ and negative if $v \in V_{\max}$.
811 Since $\Delta^{\gamma, \mathbf{r}}(a)$ is continuous in γ on $[0, 1)$ and its magnitude stays strictly positive throughout
812 $(\gamma_{\text{bw}}(\mathbf{r}), 1)$, continuity extends the strictly positive magnitude and the consistent sign to a left-
813 neighborhood $(\gamma_{\text{bw}}(\mathbf{r}) - \eta, \gamma_{\text{bw}}(\mathbf{r})]$ for some $\eta > 0$ uniform over the finite edge set. On this neigh-
814 borhood, $a^*(v)$ is the unique edge at v achieving $\Delta^{\gamma, \mathbf{r}}(a^*(v)) = 0$, so the pair $(\pi_{\min}^{\text{bw}}, \pi_{\max}^{\text{bw}})$ remains
815 optimal on $(\gamma_{\text{bw}}(\mathbf{r}) - \eta, 1)$. This contradicts the definition of $\gamma_{\text{bw}}(\mathbf{r})$ as the infimum of valid thresh-
816 olds. Therefore, $\gamma_{\text{bw}}(\mathbf{r}) \leq \max(1 - \frac{\delta}{36n^2M}, 1 - \frac{B}{10nM})$.

817 For a given $x \geq 1$, we substitute the designated $M(x) = 1 + \frac{1}{\theta} \ln(mx)$, $\delta(x) \triangleq \frac{36n^2M(x)}{x}$, and
818 $B(x) \triangleq \frac{10nM(x)}{x}$. By definition, $\frac{\delta(x)}{36n^2M(x)} = \frac{B(x)}{10nM(x)} = \frac{1}{x}$, which implies $\gamma_{\text{bw}}(\mathbf{r}) \leq \max(0, 1 - \frac{1}{x})$.
819 Under the event $\mathcal{E}_1(M(x)) \cap \mathcal{E}_2(B(x)) \cap \mathcal{E}_3(\delta(x)) \cap \mathcal{E}_4$, we have $\gamma_{\text{bw}}(\mathbf{r}) \leq 1 - \frac{1}{x}$. Therefore, the event
820 $\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}$ is upper bounded:

$$\mathbb{P}_{\xi}[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] \leq \mathbb{P}_{\xi}[\overline{\mathcal{E}_1}(M(x))] + \mathbb{P}_{\xi}[\overline{\mathcal{E}_2}(B(x))] + \mathbb{P}_{\xi}[\overline{\mathcal{E}_3}(\delta(x))] + \mathbb{P}_{\xi}[\overline{\mathcal{E}_4}].$$

821 For $\mathbb{P}_\xi[\overline{\mathcal{E}_1}(M(x))]$, we have

$$\begin{aligned}\mathbb{P}_\xi[\overline{\mathcal{E}_1}(M(x))] &\leq m \exp(-\theta(M(x) - 1)) \\ &= m \exp\left(-\theta \frac{1}{\theta} \ln(mx)\right) = \frac{1}{x},\end{aligned}$$

822 where the first inequality is due to [Eq. \(1\)](#) and a union bound over the m edges, and the first
823 equality is due to the definition of $M(x)$. For $\mathbb{P}_\xi[\overline{\mathcal{E}_2}(B(x))]$, we have

$$\begin{aligned}\mathbb{P}_\xi[\overline{\mathcal{E}_2}(B(x))] &\leq 2B(x)m^2n\phi \\ &= 2m^2n \frac{10nM(x)}{x} \phi = \frac{20m^2n^2\phi M(x)}{x},\end{aligned}$$

824 where the first inequality is due to [Lemma 5](#) and a union bound over the m edges, and the first
825 equality is due to the definition of $B(x)$. For $\mathbb{P}_\xi[\overline{\mathcal{E}_3}(\delta(x))]$, we have

$$\begin{aligned}\mathbb{P}_\xi[\overline{\mathcal{E}_3}(\delta(x))] &\leq 2m\delta(x)\phi \\ &= 2m \frac{36n^2M(x)}{x} \phi = \frac{72mn^2\phi M(x)}{x},\end{aligned}$$

826 where the first inequality is due to [Lemma 4](#) and a union bound over the m edges, and the first
827 equality is due to the definition of $\delta(x)$. We also have $\mathbb{P}_\xi[\overline{\mathcal{E}_4}] = 0$ by [Proposition 12](#). Therefore, we
828 obtain

$$\begin{aligned}\mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] &\leq \frac{1}{x} + \frac{20m^2n^2\phi M(x)}{x} + \frac{72mn^2\phi M(x)}{x} \\ &= \frac{1 + 4mn^2(5m + 18)\phi M(x)}{x} \\ &\leq \frac{1 + 92m^2n^2\phi M(x)}{x} \leq \frac{KM(x)}{x},\end{aligned}$$

829 where the second inequality uses $5m + 18 \leq 23m$ (which holds since $m \geq 1$), and the last inequality
830 uses $1 \leq 92m^2n^2\phi M(x)$ (which holds since $M(x) \geq 1$ for $x \geq 1$ and $m, n \geq 1$) and the definition of
831 K in [Eq. \(3\)](#). This completes the proof. \square

832 D.2 Restarting Policy Iteration

833 **Lemma 7** (Total iteration bound). *Let G be any game graph with $n \geq 2$, let $\gamma \in [0, 1)$, and let
834 $\mathbf{r} \in \mathbb{R}^E$. Write $N_{\text{iter}}(\mathbf{r})$ for the total number of iterations executed by $\text{RePI}(G, \gamma, \mathbf{r})$ ([Algorithm 2](#)),
835 summed over all calls to PI. Then $N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2)n^n$ for every \mathbf{r} , and, under the smoothed
836 model, with probability 1,*

$$N_{\text{iter}}(\mathbf{r}) \leq 27 \cdot \frac{m}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{e \cdot n}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right).$$

837 *Proof.* Let $k_{\text{max}} = \lceil n \ln n \rceil$ and let $k_{\text{final}} \in \{0, 1, \dots, k_{\text{max}}\}$ denote the index of the last iter-
838 ation of the for-loop executed by RePI . The algorithm performs PI at threshold γ_k for each
839 $k \in \{0, 1, \dots, k_{\text{final}}\}$ (or at the target γ when $\gamma \leq \gamma_{k_{\text{final}}}$), and optionally one fallback $\text{PI}(G, \gamma, \mathbf{r})$
840 in Phase 3. By [Proposition 1](#), each PI call at threshold $\gamma_k = 1 - e^{-k}$ executes at most
841 $\frac{6m}{1 - \gamma_k} \ln\left(\frac{n}{1 - \gamma_k}\right) = 6m e^k (\ln n + k)$ iterations, and is also trivially bounded by n^n because the number
842 of iterations is upper bounded by the number of positional strategies. Therefore, we have

$$N_{\text{iter}}(\mathbf{r}) \leq \sum_{k=0}^{k_{\text{final}}} \min(6m e^k (\ln n + k), n^n) + n^n \cdot \mathbf{1}\{\text{Phase 3 triggered}\}.$$

843 *Deterministic bound.* Bounding every term in the sum by n^n and using that there are at most
 844 $k_{\max} + 1$ terms plus at most one fallback gives $N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2) n^n$ for every \mathbf{r} .

845 *Smoothed geometric bound.* Let $\Gamma_{\text{alg}} \triangleq \{\gamma_k : k = 0, 1, \dots, k_{\max}\}$ be the finite set of thresholds used
 846 by the algorithm. By [Lemma 14](#), with probability 1, for every $\gamma_k \in \Gamma_{\text{alg}}$ with $\gamma_k > \gamma_{\text{bw}}(\mathbf{r})$, the
 847 game $G^{\gamma_k, \mathbf{r}}$ has a unique optimal strategy pair, and this pair is Blackwell optimal. Condition on
 848 this probability-one event. We bound the total in terms of $\gamma_{\text{bw}}(\mathbf{r})$, splitting into two cases according
 849 to whether the Blackwell threshold is resolved within the for-loop. In each case we establish the
 850 two bounds

$$e^{k_{\text{final}}+1} \leq \frac{e^2}{1 - \gamma_{\text{bw}}(\mathbf{r})} \quad \text{and} \quad \ln n + k_{\text{final}} \leq \ln\left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right). \quad (4)$$

851 **Case 1:** $\gamma_{\text{bw}}(\mathbf{r}) < \gamma_{k_{\max}}$. Since $\gamma_0 = 0$, there exists $i \in \{0, 1, \dots, k_{\max} - 1\}$ with $\gamma_i \leq \gamma_{\text{bw}}(\mathbf{r}) < \gamma_{i+1}$.
 852 Because $\gamma_{i+1} > \gamma_{\text{bw}}(\mathbf{r})$, the conditioned event implies that the unique optimal pair of $G^{\gamma_{i+1}, \mathbf{r}}$ is
 853 Blackwell optimal, so $\text{PI}(G, \gamma_{i+1}, \mathbf{r})$ returns a Blackwell optimal pair; being Blackwell optimal, it
 854 is optimal for every $\gamma' \in [\gamma_{\text{bw}}(\mathbf{r}), 1)$. Hence, if iteration $k = i + 1$ reaches Phase 2 (that is, if
 855 $\gamma > \gamma_{i+1}$), the returned pair is optimal for the target γ and Phase 2 succeeds; otherwise $\gamma \leq \gamma_{i+1}$
 856 and the algorithm has already returned $\text{PI}(G, \gamma, \mathbf{r})$ at some iteration $k \leq i + 1$. Either way the
 857 algorithm terminates without triggering Phase 3, so $k_{\text{final}} \leq i + 1$ and the Phase 3 term vanishes.
 858 From $1 - \gamma_i = e^{-i} \geq 1 - \gamma_{\text{bw}}(\mathbf{r})$ we obtain $e^i \leq 1/(1 - \gamma_{\text{bw}}(\mathbf{r}))$, hence

$$e^{k_{\text{final}}+1} \leq e^{i+2} \leq \frac{e^2}{1 - \gamma_{\text{bw}}(\mathbf{r})} \quad \text{and} \quad k_{\text{final}} \leq \ln \frac{1}{1 - \gamma_{\text{bw}}(\mathbf{r})} + 1,$$

859 which gives [Eq. \(4\)](#) since $\ln n + k_{\text{final}} \leq \ln n + \ln \frac{1}{1 - \gamma_{\text{bw}}(\mathbf{r})} + 1 = \ln(en/(1 - \gamma_{\text{bw}}(\mathbf{r})))$.

860 **Case 2:** $\gamma_{\text{bw}}(\mathbf{r}) \geq \gamma_{k_{\max}}$. Here $1/(1 - \gamma_{\text{bw}}(\mathbf{r})) \geq e^{k_{\max}} \geq n^n$, and the algorithm may run all $k_{\max} + 1$
 861 for-loop iterations plus a Phase 3 fallback, so $k_{\text{final}} = k_{\max}$. Hence

$$e^{k_{\text{final}}+1} = e \cdot e^{k_{\max}} \leq \frac{e}{1 - \gamma_{\text{bw}}(\mathbf{r})} \leq \frac{e^2}{1 - \gamma_{\text{bw}}(\mathbf{r})} \quad \text{and} \quad k_{\text{final}} = k_{\max} \leq \ln \frac{1}{1 - \gamma_{\text{bw}}(\mathbf{r})},$$

862 which gives [Eq. \(4\)](#) since $\ln n + k_{\text{final}} \leq \ln n + \ln \frac{1}{1 - \gamma_{\text{bw}}(\mathbf{r})} \leq \ln(en/(1 - \gamma_{\text{bw}}(\mathbf{r})))$.

863 We now bound the for-loop contribution using [Eq. \(4\)](#):

$$\begin{aligned} \sum_{k=0}^{k_{\text{final}}} 6m e^k (\ln n + k) &\leq 6m (\ln n + k_{\text{final}}) \sum_{k=0}^{k_{\text{final}}} e^k \\ &\leq 6m (\ln n + k_{\text{final}}) \cdot \frac{e^{k_{\text{final}}+1}}{e - 1} \\ &\leq \frac{6e^2}{e - 1} \cdot \frac{m}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right), \end{aligned}$$

864 where the first inequality uses $\ln n + k \leq \ln n + k_{\text{final}}$ for every $k \leq k_{\text{final}}$, the second sums the
 865 geometric series $\sum_{k=0}^{k_{\text{final}}} e^k \leq e^{k_{\text{final}}+1}/(e - 1)$, and the third substitutes both bounds of [Eq. \(4\)](#). It
 866 remains to account for the Phase 3 fallback. In Case 1 this term vanishes; in Case 2 it contributes
 867 at most

$$n^n \leq \frac{1}{1 - \gamma_{\text{bw}}(\mathbf{r})} \leq \frac{m}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right),$$

868 using $m \geq 1$ and $\ln(en) \geq 1$. Adding the two contributions yields

$$N_{\text{iter}}(\mathbf{r}) \leq \left(\frac{6e^2}{e - 1} + 1\right) \cdot \frac{m}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right) \leq 27 \cdot \frac{m}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right).$$

869 which is the second bound. As it was derived on the probability-one event from [Lemma 14](#), this
 870 bound holds with probability 1. \square

871 D.3 High Probability Guarantees

872 **Lemma 16.** *Let G be any game graph. Under the smoothed model, for any $\epsilon \in (0, 1)$, define*

$$L_\epsilon \triangleq \theta + \ln\left(1 + \frac{2mK}{\epsilon\theta}\right) \quad \text{and} \quad x_\epsilon \triangleq \frac{2KL_\epsilon}{\epsilon\theta}.$$

873 *With probability at least $1 - \epsilon$ over the perturbation \mathbf{r} , [Algorithm 2](#) executes PI for a total of at most*
 874 *$27 m x_\epsilon \ln(en x_\epsilon)$ iterations. In particular, x_ϵ is polynomial in $n, m, \phi, 1/\theta$, and $1/\epsilon$.*

875 *Proof.* Since $L_\epsilon \geq \theta$, $K \geq 1$, $m \geq 1$, and $\epsilon \in (0, 1)$, we have $x_\epsilon \geq \frac{2K}{\epsilon} > 2$.

876 We claim that $\gamma_{\text{bw}}(\mathbf{r}) \leq 1 - 1/x_\epsilon$ with probability at least $1 - \epsilon$. Indeed, applying [Lemma 6](#)
 877 with $x = x_\epsilon$ yields

$$\mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x_\epsilon}] \leq \frac{KM(x_\epsilon)}{x_\epsilon} = \frac{\epsilon}{2L_\epsilon} \left(\theta + \ln(mx_\epsilon)\right),$$

878 where the equality uses the definitions of $M(x)$ and x_ϵ . By the definition of L_ϵ ,

$$\ln(mx_\epsilon) = \ln\left(\frac{2mK}{\epsilon\theta}\right) + \ln L_\epsilon \leq \ln\left(1 + \frac{2mK}{\epsilon\theta}\right) + L_\epsilon = 2L_\epsilon - \theta,$$

879 where the inequality uses $\ln L_\epsilon \leq L_\epsilon$. Consequently,

$$\frac{KM(x_\epsilon)}{x_\epsilon} \leq \frac{\epsilon}{2L_\epsilon} (\theta + 2L_\epsilon - \theta) = \epsilon,$$

880 which yields $\gamma_{\text{bw}}(\mathbf{r}) \leq 1 - 1/x_\epsilon$ with probability at least $1 - \epsilon$. As the geometric bound of [Lemma 7](#)
 881 holds with probability 1, its intersection with this event still has probability at least $1 - \epsilon$; assume
 882 both hold for the rest of the proof.

883 The map $y \mapsto y \ln(en y)$ is increasing on $y \geq 1$. Therefore, we have

$$\frac{1}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})}\right) \leq x_\epsilon \ln(en x_\epsilon).$$

884 [Lemma 7](#) then yields $N_{\text{iter}}(\mathbf{r}) \leq 27 m x_\epsilon \ln(en x_\epsilon)$. □

885 D.4 Expected Runtime Guarantees

886 **Lemma 9** (Bounding expected truncated inverse gap). *Let G be any game graph with $n \geq 2$. Let*
 887 *K be defined in [Eq. \(3\)](#). Under the smoothed model, we have*

$$\mathbb{E}_\xi[Y(\mathbf{r})] \leq \frac{6(\theta + 1)}{\theta} K(n \ln m)^3.$$

888 *Proof.* Let x_0 be the unique real number satisfying $x_0 \geq 1$ and $x_0 \ln x_0 = n^n$. We compute the
 889 expectation of $Y(\mathbf{r})$ by integrating the tail probability and applying the substitution $y = x \ln x$,
 890 where $dy = (1 + \ln x)dx$:

$$\begin{aligned} \mathbb{E}_\xi[Y(\mathbf{r})] &= \int_0^{n^n} \mathbb{P}_\xi[Y(\mathbf{r}) > y] dy \\ &= \int_1^{x_0} \mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] (1 + \ln x) dx \\ &\leq \int_1^e 1 \cdot (1 + \ln x) dx + \int_e^{x_0} \mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] (1 + \ln x) dx \\ &= e + \int_e^{x_0} \mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] (1 + \ln x) dx. \end{aligned} \tag{5}$$

891 For the second integral, we substitute the bound $\mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}] \leq \frac{KM(x)}{x}$ established in
 892 [Lemma 6](#), which applies on $[e, x_0]$ since $x \geq e > 1$:

$$\int_e^{x_0} \mathbb{P}_\xi[\gamma_{\text{bw}}(\mathbf{r}) > 1 - \frac{1}{x}](1 + \ln x) dx \leq K \int_e^{x_0} \frac{M(x)(1 + \ln x)}{x} dx.$$

893 We apply the substitution $u = \ln x$, meaning $dx/x = du$. The integration limits shift to $u \in$
 894 $[1, \ln x_0]$. Since $x_0 \leq n^n$, the upper limit is strictly bounded by $n \ln n$. This yields:

$$\begin{aligned} \int_e^{x_0} \frac{M(x)(1 + \ln x)}{x} dx &\leq \int_1^{n \ln n} M(e^u)(1 + u) du \\ &\leq M(e^{n \ln n}) \int_1^{n \ln n} (1 + u) du \\ &= \left(1 + \frac{\ln m + n \ln n}{\theta}\right) \left[u + \frac{u^2}{2}\right]_1^{n \ln n} \\ &\leq \left(1 + \frac{\ln m + n \ln n}{\theta}\right) \left(\frac{1}{2}(n \ln n)^2 + n \ln n\right) \\ &\leq \left(1 + \frac{\ln m + n \ln n}{\theta}\right) \frac{3}{2}(n \ln n)^2 \\ &\leq \frac{\theta + 1}{\theta} (\ln m + n \ln n) \frac{3}{2}(n \ln n)^2 \\ &\leq \frac{\theta + 1}{\theta} 2n \ln m \frac{3}{2}(n \ln n)^2 \\ &\leq \frac{\theta + 1}{\theta} 2n \ln m \frac{3}{2}(n \ln m)^2 \\ &= \frac{3(\theta + 1)}{\theta} (n \ln m)^3, \end{aligned}$$

895 where the first equality follows from $M(e^u) = 1 + \frac{\ln m + u}{\theta}$; the inequality $\frac{1}{2}(n \ln n)^2 + n \ln n \leq$
 896 $\frac{3}{2}(n \ln n)^2$ uses $n \ln n \leq (n \ln n)^2$ (which holds for all $n \geq 1$, since $n \ln n \geq 0$); the bound $1 +$
 897 $\frac{\ln m + n \ln n}{\theta} \leq \frac{\theta + 1}{\theta} (\ln m + n \ln n)$ uses $1 \leq \ln m + n \ln n$ (which holds for $n \geq 2$); the bounds $\ln m +$
 898 $n \ln n \leq 2n \ln m$ and $(n \ln n)^2 \leq (n \ln m)^2$ both use $m \geq n$.

899 Substituting this bound back into [Eq. \(5\)](#) yields $\mathbb{E}_\xi[Y(\mathbf{r})] \leq e + K \frac{3(\theta + 1)}{\theta} (n \ln m)^3$. Finally, since
 900 $\frac{\theta + 1}{\theta} \geq 1$ and $n \ln m \geq 2 \ln 2 > 1$ (as $m \geq n \geq 2$), we have $\frac{3(\theta + 1)}{\theta} (n \ln m)^3 \geq 3(n \ln m)^3 > 3 > e$;
 901 since $K \geq 1$, this gives $e \leq K \frac{3(\theta + 1)}{\theta} (n \ln m)^3$, and hence $\mathbb{E}_\xi[Y(\mathbf{r})] \leq 2K \frac{3(\theta + 1)}{\theta} (n \ln m)^3$, completing
 902 the proof. \square

903 **Lemma 17.** *Let G be any game graph with $n \geq 2$. Under the smoothed model, the expected total*
 904 *number of iterations in all the PI's executed by [Algorithm 2](#) is bounded by*

$$324 m \ln n K \frac{3(\theta + 1)}{\theta} (n \ln m)^3,$$

905 where K is defined in [Eq. \(3\)](#). In particular, this bound is polynomial in n, m, ϕ , and $1/\theta$.

906 *Proof.* Write $x \triangleq 1/(1 - \gamma_{\text{bw}}(\mathbf{r})) \geq 1$, and recall the truncated inverse gap $Y(\mathbf{r}) = \min(x \ln x, n^n)$.
 907 By [Lemma 7](#), with probability 1 we have

$$N_{\text{iter}}(\mathbf{r}) \leq \min\left(27 m x \ln(enx), (\lceil n \ln n \rceil + 2) n^n\right).$$

908 We claim that

$$N_{\text{iter}}(\mathbf{r}) \leq 27 m ((\ln n + 2) Y(\mathbf{r}) + e(\ln n + 1)). \quad (6)$$

909 Indeed, we prove this by splitting into two cases according to the truncation.

910 **Case 1:** $x \ln x \leq n^n$. Here $Y(\mathbf{r}) = x \ln x$, and we bound the minimum by its first term:

$$\begin{aligned}
N_{\text{iter}}(\mathbf{r}) &\leq 27 m x \ln(en x) \\
&= 27 m(x \ln x + (\ln n + 1)x) \\
&\leq 27 m(Y(\mathbf{r}) + (\ln n + 1)(Y(\mathbf{r}) + e)) \\
&= 27 m((\ln n + 2)Y(\mathbf{r}) + e(\ln n + 1)),
\end{aligned}$$

911 where the first equality uses $\ln(en x) = \ln x + \ln n + 1$, and the second inequality uses $x \ln x = Y(\mathbf{r})$
912 together with $x \leq x \ln x + e = Y(\mathbf{r}) + e$, valid for all $x \geq 1$.

913 **Case 2:** $x \ln x > n^n$. Here $Y(\mathbf{r}) = n^n$, and we bound the minimum by its second term. Since
914 every vertex has at least one outgoing edge, $m \geq n$, hence $\lceil n \ln n \rceil + 2 \leq n \ln n + 3 \leq 27 m(\ln n + 2)$.
915 Therefore

$$N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2) n^n \leq 27 m(\ln n + 2) n^n \leq 27 m((\ln n + 2)Y(\mathbf{r}) + e(\ln n + 1)).$$

916 Consequently, [Eq. \(6\)](#) holds with probability 1; since $N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2) n^n$ for every \mathbf{r} by
917 [Lemma 7](#), the complementary null event contributes nothing to $\mathbb{E}_\xi[N_{\text{iter}}(\mathbf{r})]$. Taking the expectation
918 of [Eq. \(6\)](#) over the smoothed perturbation \mathbf{r} and applying [Lemma 9](#),

$$\begin{aligned}
\mathbb{E}_\xi[N_{\text{iter}}(\mathbf{r})] &\leq 27 m((\ln n + 2) \mathbb{E}_\xi[Y(\mathbf{r})] + e(\ln n + 1)) \\
&\leq 27 m\left(2(\ln n + 2)K \frac{3(\theta + 1)}{\theta} (n \ln m)^3 + e(\ln n + 1)\right) \\
&\leq 81 m(\ln n + 2)K \frac{3(\theta + 1)}{\theta} (n \ln m)^3 \\
&\leq 324 m \ln n K \frac{3(\theta + 1)}{\theta} (n \ln m)^3,
\end{aligned}$$

919 where the third inequality uses $e(\ln n + 1) \leq (\ln n + 2)K \frac{3(\theta + 1)}{\theta} (n \ln m)^3$, which follows from $e \leq$
920 $K \frac{3(\theta + 1)}{\theta} (n \ln m)^3$ (since $K \geq 1$ and $\frac{3(\theta + 1)}{\theta} (n \ln m)^3 \geq 3(n \ln m)^3 \geq 3 > e$), and the last inequality
921 uses $\ln n + 2 \leq 4 \ln n$ (which holds for $n \geq 2$). This is the claimed bound. \square

922 E Mean-payoff Games

923 In this section, we provide the technical development for mean-payoff games that establishes [Corol-](#)
924 [lary 11](#), combining a high-probability guarantee ([Lemma 19](#)) and an expected runtime guarantee
925 ([Lemma 20](#)). We first adapt the restarting algorithm of [Section 7.2](#) to the mean-payoff setting and
926 establish a total iteration bound for it ([Section E.1](#)). Combining this bound with the tail bound
927 of [Section 7.1](#) yields smoothed polynomial runtime with high probability ([Section E.2](#)); combining
928 it with the expected truncated inverse gap of [Section 7.4](#) yields smoothed polynomial expected
929 runtime ([Section E.3](#)).

930 E.1 Restarting Mean-payoff Policy Iteration

931 In this subsection, we adapt the restarting variant of policy iteration to the mean-payoff setting,
932 called RePI_{mp} , presented in full detail in [Algorithm 3](#), and then establish a deterministic bound on
933 its total iteration count in terms of the Blackwell threshold $\gamma_{\text{bw}}(\mathbf{r})$ ([Lemma 18](#)).

934 **Informal description.** Our algorithm sweeps through the same geometric sequence of thresholds
935 $\gamma_k = 1 - e^{-k}$ for $k = 0, 1, \dots, \lceil n \ln n \rceil$, with each successive threshold e times closer to 1 than the
936 previous one. For each k in turn, it solves the game with the discount factor γ_k via the standard
937 discounted-sum $\text{PI}(G, \gamma_k, \mathbf{r})$ to obtain a candidate strategy pair (π_{\min}, π_{\max}) . The algorithm then

938 tests whether this candidate pair is already optimal for the mean-payoff game. For this purpose,
 939 we use the combinatorial strategy evaluation of [BV07], which reduces mean-payoff games to the
 940 longest shortest paths problem and evaluates strategies using shortest distances to a sink. Their
 941 method can determine whether a given positional strategy pair is mean-payoff optimal; we invoke
 942 this test as a black box. If the test succeeds, the algorithm computes the mean-payoff value vector
 943 induced by the candidate pair and returns both. Otherwise it proceeds to the next threshold γ_{k+1}
 944 and retries. If every threshold in the sequence fails to be optimal, the algorithm falls back to the
 945 longest-shortest-path strategy-improvement algorithm of [BV07], denoted by $\text{Pl}_{\text{mp}}(G, \mathbf{r})$. We use
 946 this algorithm as a deterministic black-box fallback: it returns a mean-payoff optimal positional
 947 strategy pair and value vector, and for our purposes the upper bound of at most n^n strategy
 948 improvements is sufficient, which is due to the fact that there are at most n^n positional strategies
 949 pairs in a game with n vertices.

Algorithm 3 Restarting Mean-payoff Policy Iteration: $\text{RePl}_{\text{mp}}(G, \mathbf{r})$

```

1: Input: Graph  $G = (V_{\min}, V_{\max}, E)$  with  $n = |V| \geq 2$ , reward vector  $\mathbf{r} \in \mathbb{R}^E$ .
2: Output: Mean-payoff optimal strategy pair  $(\pi_{\min}^*, \pi_{\max}^*)$  and mean-payoff value vector  $\mu^{\text{mp}, \mathbf{r}}$ 
   for the game  $G^{\text{mp}, \mathbf{r}}$ .
3: Let  $k_{\max} = \lceil n \ln n \rceil$ .
4: for  $k = 0, 1, \dots, k_{\max}$  do
5:   Set the threshold discount factor  $\gamma_k = 1 - e^{-k}$ .
6:   // Phase 1: Solve the discounted game at the threshold
7:    $(\pi_{\min}, \pi_{\max}, \mu_{\gamma_k}) \leftarrow \text{Pl}(G, \gamma_k, \mathbf{r})$ 
8:   // Phase 2: Mean-payoff optimality test
9:   Use the longest shortest paths algorithm of [BV07] to test whether  $(\pi_{\min}, \pi_{\max})$  is mean-
   payoff optimal.
10:  if  $(\pi_{\min}, \pi_{\max})$  is mean-payoff optimal then
11:    Compute the mean-payoff value vector  $\mu^{\text{mp}, \mathbf{r}}$  induced by  $(\pi_{\min}, \pi_{\max})$ .
12:    return  $(\pi_{\min}, \pi_{\max}, \mu^{\text{mp}, \mathbf{r}})$ 
13:  end if
14: end for
15: // Phase 3: Fallback
16: return  $\text{Pl}_{\text{mp}}(G, \mathbf{r})$ 

```

950 The intuition mirrors the discounted-sum case. Under the smoothed model, with high prob-
 951 ability some threshold γ_k in the sequence strictly exceeds the Blackwell threshold $\gamma_{\text{bw}}(\mathbf{r})$. Thus,
 952 the strategy pair returned by $\text{Pl}(G, \gamma_k, \mathbf{r})$ is Blackwell optimal and consequently mean-payoff op-
 953 timal [BK76], so the Phase 2 test succeeds and the algorithm terminates without reaching the
 954 fallback.

955 **Lemma 18** (Total iteration bound). *Let G be any game graph with $n \geq 2$ and let $\mathbf{r} \in \mathbb{R}^E$; write*
 956 *$N_{\text{iter}}(\mathbf{r})$ for the total number of iterations performed by $\text{RePl}_{\text{mp}}(G, \mathbf{r})$ (Algorithm 3), across all calls*
 957 *to Pl and the optional Pl_{mp} fallback. Then $N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2) n^n$ for every \mathbf{r} , and, under the*
 958 *smoothed model, with probability 1,*

$$N_{\text{iter}}(\mathbf{r}) \leq 27 \cdot \frac{m}{1 - \gamma_{\text{bw}}(\mathbf{r})} \ln \left(\frac{en}{1 - \gamma_{\text{bw}}(\mathbf{r})} \right).$$

959 *Proof.* The deterministic bound $N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2) n^n$ holds exactly as in Lemma 7: there are
 960 at most $\lceil n \ln n \rceil + 2$ calls, each terminating in at most n^n iterations (the fallback $\text{Pl}_{\text{mp}}(G, \mathbf{r})$ likewise,
 961 by [BV07]). For the geometric bound, condition on the probability-one event from Lemma 14 with

962 $\Gamma_{\text{alg}} = \{\gamma_k : k = 0, 1, \dots, k_{\text{max}}\}$, and run the same case analysis as in [Lemma 7](#) with two minor
963 changes. First, when $\gamma_{i+1} > \gamma_{\text{bw}}(\mathbf{r})$, the unique optimal pair of $G^{\gamma_{i+1}, \mathbf{r}}$ is Blackwell optimal and
964 is also *mean-payoff* optimal [[BK76](#)], so $\text{PI}(G, \gamma_{i+1}, \mathbf{r})$ returns it; the test based on longest shortest
965 paths from [[BV07](#)] therefore certifies the pair as mean-payoff optimal in Phase 2, and the algorithm
966 terminates without triggering Phase 3. Second, when Phase 3 is triggered (i.e., $\gamma_{\text{bw}}(\mathbf{r}) \geq \gamma_{k_{\text{max}}}$), the
967 fallback is $\text{PI}_{\text{mp}}(G, \mathbf{r})$ rather than $\text{PI}(G, \gamma, \mathbf{r})$; it terminates in at most n^n iterations [[BV07](#)], and the
968 same bound $n^n \leq \frac{m}{1-\gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1-\gamma_{\text{bw}}(\mathbf{r})}\right)$ applies in this regime. The geometric sum then combines
969 as in [Lemma 7](#) to give the second bound, which holds with probability 1 by [Lemma 14](#). \square

970 E.2 High Probability Guarantees

971 In this subsection, we combine the total iteration bound of [Lemma 18](#) with the tail bound from
972 [Section 7.1](#) to obtain a high-probability bound on the number of iterations of RePI_{mp} .

973 **Lemma 19.** *Let G be any game graph. Under the smoothed model, for any $\epsilon \in (0, 1)$, let x_ϵ be as*
974 *defined in [Lemma 16](#). With probability at least $1 - \epsilon$ over the perturbation \mathbf{r} , RePI_{mp} ([Algorithm 3](#))*
975 *executes a total of at most $27 m x_\epsilon \ln(en x_\epsilon)$ iterations.*

976 *Proof.* By Step 1 of the proof of [Lemma 16](#), $\gamma_{\text{bw}}(\mathbf{r}) \leq 1 - 1/x_\epsilon$ with probability at least $1 - \epsilon$. As the
977 geometric bound of [Lemma 18](#) holds with probability 1, assume this event together with it (their
978 intersection still has probability at least $1 - \epsilon$), so $\frac{1}{1-\gamma_{\text{bw}}(\mathbf{r})} \leq x_\epsilon$, and the map $y \mapsto y \ln(en y)$ being
979 increasing on $y \geq 1$ gives $\frac{1}{1-\gamma_{\text{bw}}(\mathbf{r})} \ln\left(\frac{en}{1-\gamma_{\text{bw}}(\mathbf{r})}\right) \leq x_\epsilon \ln(en x_\epsilon)$. [Lemma 18](#) then yields $N_{\text{iter}}(\mathbf{r}) \leq$
980 $27 m x_\epsilon \ln(en x_\epsilon)$. \square

981 E.3 Expected Runtime Guarantees

982 In this subsection, we combine the total iteration bound of [Lemma 18](#) with the bound on the
983 expected truncated inverse gap from [Lemma 9](#) to obtain a bound on the expected number of
984 iterations of RePI_{mp} .

985 **Lemma 20.** *Let G be any game graph. Under the smoothed model, the expected total number of*
986 *iterations executed by [Algorithm 3](#) (across all PI executions and the fallback PI_{mp} , when triggered)*
987 *is bounded by*

$$324 m \ln n K \frac{3(\theta + 1)}{\theta} (n \ln m)^3,$$

988 *where K is defined in [Eq. \(3\)](#). In particular, this bound is polynomial in n, m, ϕ , and $1/\theta$.*

989 *Proof.* The argument is identical to the proof of [Lemma 17](#), invoking the total iteration bound
990 of [Lemma 18](#) in place of [Lemma 7](#). With probability 1 the same pointwise bound $N_{\text{iter}}(\mathbf{r}) \leq$
991 $27 m ((\ln n + 2)Y(\mathbf{r}) + e(\ln n + 1))$ holds, and since $N_{\text{iter}}(\mathbf{r}) \leq (\lceil n \ln n \rceil + 2)n^n$ for every \mathbf{r} the com-
992plementary null event contributes nothing to the expectation; taking the expectation and applying
993 [Lemma 9](#) yields the stated bound. \square

994 *Proof of [Corollary 11](#).* The high-probability guarantee follows from [Lemma 19](#). The expected run-
995time guarantee follows from [Lemma 20](#). \square