

# Sample-based Approximation of Nash in Large Many-Player Games via Gradient Descent

Ian Gemp  
DeepMind  
London, United Kingdom  
imgemp@deepmind.com

Yoram Bachrach  
DeepMind  
London, United Kingdom  
yorambac@deepmind.com

Andrea Tacchetti  
DeepMind  
London, United Kingdom  
atacchet@deepmind.com

Rahul Savani  
University of Liverpool  
Liverpool, United Kingdom  
rahul.savani@liverpool.ac.uk

Thomas Anthony  
DeepMind  
London, United Kingdom  
twa@deepmind.com

Tom Eccles  
DeepMind  
London, United Kingdom  
eccles@deepmind.com

Marc Lanctot  
DeepMind  
Edmonton, Canada  
lanctot@deepmind.com

Richard Everett  
DeepMind  
London, United Kingdom  
reverett@deepmind.com

János Kramár  
DeepMind  
London, United Kingdom  
janosk@deepmind.com

## ABSTRACT

Nash equilibrium is a central concept in game theory. Several Nash solvers exist, yet none scale to normal-form games with many actions and many players, especially those with payoff tensors too big to be stored in memory. In this work, we propose an approach that iteratively improves an approximation to a Nash equilibrium through joint play. It accomplishes this by tracing a previously established homotopy that defines a continuum of equilibria for the game regularized with decaying levels of entropy. This continuum asymptotically approaches the *limiting logit equilibrium*, proven by McKelvey and Palfrey (1995) to be unique in *almost* all games, thereby partially circumventing the well-known equilibrium selection problem of many-player games. To encourage iterates to remain near this path, we efficiently minimize *average deviation incentive* via stochastic gradient descent, intelligently sampling entries in the payoff tensor as needed. Monte Carlo estimates of the stochastic gradient from joint play are biased due to the appearance of a nonlinear max operator in the objective, so we introduce additional innovations to the algorithm to alleviate gradient bias. The descent process can also be viewed as repeatedly constructing and reacting to a polymatrix approximation to the game. In these ways, our proposed approach, *average deviation incentive descent with adaptive sampling* (ADIDAS), is most similar to three classical approaches, namely homotopy-type, Lyapunov, and iterative polymatrix solvers. The lack of local convergence guarantees for biased gradient descent prevents guaranteed convergence to Nash, however, we demonstrate through extensive experiments the ability of this approach to approximate a unique Nash equilibrium in normal-form games with as many as seven players and twenty one actions (several billion outcomes) that are orders of magnitude larger than those possible with prior algorithms.

## KEYWORDS

Nash; Quantal Response Equilibrium; Limiting Logit Equilibrium; Homotopy; N-player; Normal-form; Empirical Game Theory

## ACM Reference Format:

Ian Gemp, Rahul Savani, Marc Lanctot, Yoram Bachrach, Thomas Anthony, Richard Everett, Andrea Tacchetti, Tom Eccles, and János Kramár. 2022. Sample-based Approximation of Nash in Large Many-Player Games via Gradient Descent. In *Proc. of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2022), Online, May 9–13, 2022*, IFAA-MAS, 30 pages.

## 1 INTRODUCTION

Core concepts from game theory underpin many advances in multi-agent systems research. Among these, Nash equilibrium is particularly prevalent. Despite the difficulty of computing a Nash equilibrium [20, 21], a plethora of algorithms [12, 33, 41, 53, 54] and suitable benchmarks [50] have been developed, however, none address large normal-form games with many actions and many players, especially those too big to be stored in memory.

In this work, we develop an algorithm for approximating a Nash equilibrium of a normal-form game with so many actions and players that only a small subset of the possible outcomes in the game can be accessed at a time. We refer the reader to McKelvey and McLennan [43] for a review of approaches for normal-form games. Several algorithms exactly compute a Nash equilibrium for small normal-form games and others efficiently approximate Nash equilibria for special game classes, however, practical algorithms for approximating Nash in large normal-form games with many players, e.g. 7, and many actions, e.g., 21, is lacking. Computational efficiency is of paramount importance for large games because a general normal-form game with  $n$  players and  $m$  actions contains  $nm^n$  payoffs; simply enumerating all payoffs can be intractable and renders classical approaches ineligious. A common approach is to return the profile found by efficient no-regret algorithms that sample payoffs as needed [11] although Flokas et al. [30] recently

proved that many from this family do not converge to mixed Nash equilibria in *all* games, 2-player games included.

While significant progress has been made for computing Nash in 2-player normal-form games which can be represented as a *linear* complementarity problem, the many-player setting induces a *nonlinear* complementarity problem, which is “often hopelessly impractical to solve exactly” ([55], p. 105).<sup>1</sup> The combination of high dimensionality ( $m^n$  vs  $m^2$  distinct outcomes) and nonlinearity (utilities are degree- $n$  polynomials in the strategies vs degree-2) makes many-player games much more complex.

This more general problem arises in cutting-edge multiagent research when learning [35] and evaluating [2] agents in Diplomacy, a complex 7-player board game. Gray et al. [35] used no-regret learning to approximate a Nash equilibrium of subsampled games, however, this approach is brittle as we show later in Figure 4. In [2], five Diplomacy bots were ranked according to their mass under an approximate Nash equilibrium. We extend that work to encourage convergence to a particular Nash and introduce sampling along with several technical contributions to scale evaluation to 21 Diplomacy bots, a >1000-fold increase in meta-game size.

Equilibrium computation has been an important component of AI in multi-agent systems [55]. It has been (and remains) a critical component of super-human AI in poker [14–16, 49]. As mentioned above, Nash computation also arises when strategically summarizing a larger domain by learning a lower dimensionality empirical game [62]; such an approach was used in the AlphaStar League, leading to an agent that beat humans in StarCraft [59, 60]. Ultimately, this required solving for the Nash of a 2-player, 888-action game, which can take several seconds using state-of-the-art solvers on modern hardware. In contrast, solving an empirical game of Diplomacy, e.g., a 7-player 888-action game, would naively take longer than the current age of the universe. This is well beyond the size of any game we inspect here, however, we approximate the Nash of games several orders of magnitude larger than previously possible, thus taking a step towards this ambitious goal.

**Our Contribution:** We introduce stochastic optimization into a classical *homotopy* approach resulting in an algorithm that avoids the need to work with the full payoff tensor all at once and is, to our knowledge, the first algorithm generally capable of practically approximating a unique Nash equilibrium in large (billions of outcomes) many-player, many-action normal-form games. We demonstrate our algorithm on 2, 3, 4, 6, 7 and 10 player games (10 in Appx. I; others in §5). We also perform various ablation studies of our algorithm (Appx. F), compare against several baselines including solvers from the popular Gambit library (more in Appx. H), and examine a range of domains (more in Appx. I).

The paper is organized as follows. After formulating the Nash equilibrium problem for a general  $n$ -player normal-form game, we review previous work. We discuss how we combine the insights of classical algorithms with ideas from stochastic optimization to develop our final algorithm, *average deviation incentive descent with adaptive sampling*, or ADIDAS. Finally, we compare our proposed algorithm against previous approaches on large games of interest from the literature: games such as Colonel Blotto [3], classical Nash

<sup>1</sup>While any  $n$ -player game can, in theory, be efficiently solved for approximate equilibria by reducing it to a two-player game, in practice this approach is not feasible for solving large games due to the blowups involved in the reductions. Details in Appx. B.

benchmarks from the GAMUT library [50], and games relevant to recent success on the 7-player game Diplomacy [2, 35].

## 2 PRELIMINARIES

In a finite  $n$ -player game in normal form, each player  $i \in \{1, \dots, n\}$  is given a strategy set  $\mathcal{A}_i = \{a_{i1}, \dots, a_{im_i}\}$  consisting of  $m_i$  pure strategies. The pure strategies can be naturally indexed by non-negative integers, so we redefine  $\mathcal{A}_i = \{0, \dots, m_i - 1\}$  as an abuse of notation for convenience. Each player  $i$  is also given a payoff or utility function,  $u_i : \mathcal{A} \rightarrow \mathbb{R}$  where  $\mathcal{A} = \prod_i \mathcal{A}_i$ . In games where the cardinality of each player’s strategy set is the same, we drop the subscript on  $m_i$ . Player  $i$  may play a mixed strategy by sampling from a distribution over their pure strategies. Let player  $i$ ’s mixed strategy be represented by a vector  $x_i \in \Delta^{m_i-1}$  where  $\Delta^{m_i-1}$  is the  $(m_i - 1)$ -dimensional probability simplex embedded in  $\mathbb{R}^{m_i}$ . Each function  $u_i$  is then extended to this domain so that  $u_i(\mathbf{x}) = \sum_{\mathbf{a} \in \mathcal{A}} u_i(\mathbf{a}) \prod_j x_{ja_j}$  where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $a_j \in \mathcal{A}_j$  denotes player  $j$ ’s component of the joint action  $\mathbf{a} \in \mathcal{A}$ . For convenience, let  $x_{-i}$  denote all components of  $\mathbf{x}$  belonging to players other than player  $i$ .

We say  $\mathbf{x} \in \prod_i \Delta^{m_i-1}$  is a Nash equilibrium iff, for all  $i \in \{1, \dots, n\}$ ,  $u_i(z_i, x_{-i}) \leq u_i(\mathbf{x})$  for all  $z_i \in \Delta^{m_i-1}$ , i.e., no player has any incentive to unilaterally deviate from  $\mathbf{x}$ . Nash is most commonly relaxed with  $\epsilon$ -Nash, an additive approximation:  $u_i(z_i, x_{-i}) \leq u_i(\mathbf{x}) + \epsilon$  for all  $z_i \in \Delta^{m_i-1}$ . Later we explore the idea of regularizing utilities with a function  $S_i^\tau$  (e.g., entropy) as follows:

$$u_i^\tau(\mathbf{x}) = u_i(\mathbf{x}) + S_i^\tau(x_i, x_{-i}). \quad (1)$$

As an abuse of notation, let the atomic action  $a_i$  also denote the  $m_i$ -dimensional “one-hot” vector with all zeros aside from a 1 at index  $a_i$ ; its use should be clear from the context. And for convenience, denote by  $H_{il}^i = \mathbb{E}_{x_{-il}} [u_i(a_i, a_l, x_{-il})]$  the Jacobian<sup>2</sup> of player  $i$ ’s utility with respect to  $x_i$  and  $x_l$ ;  $x_{-il}$  denotes all strategies belonging to players other than  $i$  and  $l$  and  $u_i(a_i, a_l, x_{-il})$  separates out  $l$ ’s strategy  $x_l$  from the rest of the players  $x_{-i}$ . We also introduce  $\nabla_{x_i}^l$  as player  $i$ ’s utility gradient. Note player  $i$ ’s utility can now be written succinctly as  $u_i(x_i, x_{-i}) = x_i^\top \nabla_{x_i}^i = x_i^\top H_{il}^i x_l$  for any  $l$ .

In a polymatrix game, interactions between players are limited to local, pairwise games, each of which is represented by matrices  $H_{ij}^i$  and  $H_{ij}^j$ . This reduces the exponential  $nm^n$  payoffs required to represent a general normal form game to a quadratic  $n(n-1)m^2$ , an efficiency we leverage later.

### 2.1 Related work

Several approaches exist for computing Nash equilibria of  $n$ -player normal form games<sup>3</sup>. Simplicial Subdivision (SD) [58] searches for an equilibrium over discretized simplices; accuracy depends on the grid size which scales exponentially with the number of player actions. Govindan and Wilson [33] propose a homotopy method (GW) that begins with the unique Nash distribution of a game whose payoff tensor has been perturbed by an arbitrary constant tensor. GW then scales back this perturbation while updating the Nash

<sup>2</sup>See Appx. E.1 for an example derivation of the gradient if this form is unfamiliar.

<sup>3</sup>Note that Double-Oracle [46] and PSRO [40] can be extended to  $n$ -player games, but require an  $n$ -player normal form meta-solver (Nash-solver) and so cannot be considered solvers in their own right. This work provides an approximate meta-solver.

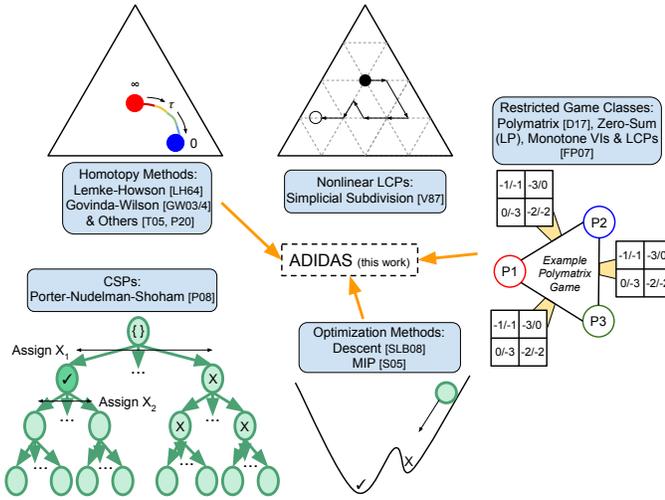


Figure 1: Algorithm Comparison and Overview.

to that of the transformed game. GW is considered an extension of the classic Lemke-Howson algorithm (1964) to 3+ player games (see §4.3, p. 107 of [55]). Another homotopy approach perturbs the payoffs with entropy bonuses, and evolves the Nash distribution along a continuum of quantal response equilibria (QREs) using a predictor-corrector method to integrate a differential equation [57] —we also aim to follow this same continuum. In a slightly different approach, Perolat et al. [52] propose an adaptive regularization scheme that repeatedly solves for the equilibrium of a transformed game. Simple search methods [53] that approach Nash computation as a constraint satisfaction problem appear to scale better than GW and SD as measured on GAMUT benchmarks [50]. Lyapunov approaches minimize non-convex energy functions with the property that zero energy implies Nash [55], however these approaches may suffer from convergence to local minima with positive energy. In some settings, such as polymatrix games with payoffs in  $[0, 1]$ , gradient descent on appropriate energy functions<sup>4</sup> guarantees a  $(\frac{1}{2} + \delta)$ -Nash in time polynomial in  $\frac{1}{\delta}$  [23] and performs well in practice [22].

*Teaser.* Our proposed algorithm consists of two key conceptual schemes. One lies at the crux of homotopy methods (see Figures 1 and 2). We initialize the Nash approximation,  $\mathbf{x}$ , to the joint uniform distribution, the unique Nash of a game with infinite-temperature entropy regularization. The temperature is then annealed over time. To recover the Nash at each temperature, we minimize an appropriately adapted energy function via (biased) stochastic gradient descent. This minimization approach can be seen as simultaneously learning a suitable polymatrix decomposition of the game similarly to Govindan and Wilson [34] but from batches of stochastic play, i.e., we compute Monte Carlo estimates of the payoffs in the bimatrix game between every pair of players by observing the outcomes of the players’ joint actions (sampled from  $\mathbf{x}$  after each update) rather than computing payoffs as exact expectations.

<sup>4</sup>Equation (2) but with max instead of  $\sum$  over player regrets. Note that for symmetric games with symmetric equilibria, these are equivalent up to a multiplicative factor  $n$ .

### 3 DEVIATION INCENTIVE & WARM-UP

We propose minimizing the energy function in equation (2) below, *average deviation incentive* (ADI), to approximate a Nash equilibrium of a large, entropy-regularized normal form game. This loss measures, on average, how much a single agent can exploit the rest of the population by deviating from a joint strategy  $\mathbf{x}$ . For sake of exposition, we drop the normalizing constant from the denominator (number of players,  $n$ ), and consider the sum instead of the average. This quantity functions as a *loss* that can be minimized over  $\mathcal{X} = \prod_i \Delta^{m_i-1}$  to find a Nash distribution. Note that when ADI is zero,  $\mathbf{x}$  is a Nash. Also, if  $\sum_k$  is replaced by  $\max_k$ , this loss measures the  $\epsilon$  of an  $\epsilon$ -Nash, and therefore, equation (2) is an upper bound on this  $\epsilon$ . Lastly, note that, in general, this loss function is non-convex and so convergence to local, suboptimal minima is theoretically possible if naively minimizing via first order methods like gradient descent —we explain in §3.1 how we circumvent this pitfall via temperature annealing. Let  $BR_k = BR(\mathbf{x}_{-k}) = \arg \max_{z_k \in \Delta^{m_k-1}} u_k^\tau(z_k, \mathbf{x}_{-k})$  be player  $k$ ’s best response to all other players’ current strategies where  $u_k^\tau$  is player  $k$ ’s utility regularized by entropy with temperature  $\tau$  and formally define

$$\mathcal{L}_{adi}^\tau(\mathbf{x}) = \sum_k \overbrace{u_k^\tau(BR_k, \mathbf{x}_{-k}) - u_k^\tau(\mathbf{x}_k, \mathbf{x}_{-k})}^{\text{incentive to deviate to } BR_k \text{ vs } \mathbf{x}_k}. \quad (2)$$

If  $\tau = 0$ , we drop the superscript and use  $\mathcal{L}_{adi}$ . The Nash equilibrium of the game regularized with Shannon entropy is called a *quantal response equilibrium*, QRE( $\tau$ ) (see p. 152-154, 343 of [31]).

Average deviation incentive has been interpreted as a pseudo-distance from Nash in prior work, where it is referred to as NashConv [40]. We prefer average deviation incentive because it more precisely describes the function and allows room for exploring alternative losses in future work. The objective can be decomposed into terms that depend on  $x_k$  (second term) and  $\mathbf{x}_{-k}$  (both terms). Minimizing the second term w.r.t.  $x_k$  seeks strategies with high utility, while minimizing both terms w.r.t.  $\mathbf{x}_{-k}$  seeks strategies that cannot be exploited by player  $k$ . In reducing  $\mathcal{L}_{adi}$ , each player  $k$  seeks a strategy that not only increases their payoff but also removes others’ temptation to exploit them.

A related algorithm is Exploitability Descent (ED) [42]. Rather than minimizing  $\mathcal{L}_{adi}$ , each player independently maximizes their utility assuming the other players play their best responses. In the two-player normal-form setting, ED is equivalent to extra-gradient [38] (see Appx. K.2). However, ED is only guaranteed to converge to Nash in two-player, zero-sum games. We include a comparison against ED as well as Fictitious-play, another popular multiagent algorithm, in Appx. H.1. We also relate  $\mathcal{L}_{adi}$  to Consensus optimization [48] in Appx. K.1.

#### 3.1 Warm-Up

McKelvey and Palfrey [45] proved the existence of a continuum of QREs starting at the uniform distribution (infinite temperature) and ending at what they called the *limiting logit equilibrium* (LLE). Furthermore, they showed this path is unique for *almost all games*,

partially circumventing the equilibrium selection problem. We encourage the reader to look ahead at Figure 2 for a visual of the homotopy that may prove helpful for the ensuing discussions.

In this work, we assume we are given one of these common games with a unique path (no branching points) so that the LLE is well defined (**Assumption 1**). Furthermore, we assume there exist no “turning points” in the temperature  $\tau$  along the continuum (**Assumption 2**). Turocy [57] explains that even in generic games, temperature might have to be temporarily increased in order to remain on the path (principal branch) to the LLE. However, Turocy also proves there exists a  $\tau^*$  such that no turning points exist with  $\tau > \tau^*$  suggesting that as long as we remain near the principal branch after  $\tau^*$ , we can expect to proceed to the LLE.

We follow the principal path by alternating between annealing the temperature and re-solving for the Nash at that temperature by minimizing  $\mathcal{L}_{adi}^\tau$ . We present a basic version of our approach that converges to the limiting logit equilibrium assuming access to exact gradients in Algorithm 1 (proof in Appx. D). We substitute  $\lambda = \frac{1}{\tau}$  and initialize  $\lambda = 0$  in order to begin at infinite temperature. The proof of this simple warm-up algorithm relies on the detailed examination of the continuum of QREs proposed in [45] and further analyzed in [57]. Theorem 1 presented below is essentially a succinct repetition of one of their known results (Assumptions 3 and 4 below are expanded on in Appx. D). In subsequent sections, we relax the exact gradient assumption and assume gradients are estimated from stochastic play (i.e., each agent samples an action from their side of the current approximation to the Nash).

---

**Algorithm 1** Warm-up: Anneal & Descend

---

- 1: Given: Total anneal steps  $T_\lambda$ , total optimizer iterations  $T^*$ , and anneal step size  $\Delta\lambda$ .
  - 2:  $\lambda = 0$
  - 3:  $\mathbf{x} \leftarrow \{\frac{1}{m_i} \mathbf{1} \forall i\}$
  - 4: **for**  $t_\lambda = 1 : T_\lambda$  **do**
  - 5:    $\lambda \leftarrow \lambda + \Delta\lambda$
  - 6:    $\mathbf{x} \leftarrow \text{OPT}(\text{loss} = \mathcal{L}_{adi}^{\tau=\lambda^{-1}}, \mathbf{x}_{init} = \mathbf{x}, \text{iters} = T^*)$
  - 7: **end for**
  - 8: **return**  $\mathbf{x}$
- 

**THEOREM 1.** *Make assumptions 1 and 2. Also, assume the QREs along the homotopy path have bounded sensitivity to  $\lambda$  given by a parameter  $\sigma$  (Assumption 3), and basins of attraction with radii lower bounded by  $r$  (Assumption 4). Let the step size  $\Delta\lambda \leq \sigma(r - \epsilon)$  with tolerance  $\epsilon$ . And let  $T^*$  be the supremum over all  $T$  such that Assumption 4 is satisfied for any inverse temperature  $\lambda \geq \Delta\lambda$ . Then, assuming gradient descent for OPT, Algorithm 1 converges to the limiting logit equilibrium  $\mathbf{x}_{\lambda=\infty}^* = \mathbf{x}_{\tau=0}^*$  in the limit as  $T_\lambda \rightarrow \infty$ .*

### 3.2 Evaluating $\mathcal{L}_{adi}^\tau$ with Joint Play

In the warm up, we assumed we could compute exact gradients which required access to the entire payoff tensor. However, we want to solve very large games where enumerating the payoff tensor is prohibitively expensive. Therefore, we are particularly interested in minimizing  $\mathcal{L}_{adi}^\tau$  when only given access to samples of joint play,  $\mathbf{a} \sim \prod_i x_i$ . The best response operator, BR, is nonlinear and

hence can introduce bias if applied to random samples. For example, consider the game given in Table 1 and assume  $x_2 = [0.5, 0.5]^\top$ .

$u_1$	$a_{21}$	$a_{22}$	$u_2$	$a_{21}$	$a_{22}$
$a_{11}$	0	0	$a_{11}$	0	0
$a_{12}$	1	-2	$a_{12}$	0	0
$a_{13}$	-2	1	$a_{13}$	0	0

**Table 1: A 2-player game with biased stochastic BR’s.**

Consider computing (row) player 1’s best response to a single action sampled from (column) player 2’s strategy  $x_2$ . Either  $a_{21}$  or  $a_{22}$  will be sampled with equal probability, which results in a best response of either  $a_{12}$  or  $a_{13}$  respectively. However, the true expected utilities for each of player 1’s actions given player 2’s strategy are  $[0, -0.5, -0.5]$  for which the best response is the first index,  $a_{11}$ . The best response operator completely filters out information on the utility of the true best response  $a_{11}$ . Intuitively, a *soft* best response operator, demonstrated in equations (3)-(5), that allows some utility information for each of the actions to pass through could alleviate the problem:

$$\mathbb{E}[\text{BR}^{\tau \rightarrow 0}] = [0.00, 0.50, 0.50] \quad (3)$$

$$\mathbb{E}[\text{BR}^{\tau=1}] \approx [0.26, 0.37, 0.37] \quad (4)$$

$$\mathbb{E}[\text{BR}^{\tau=10}] \approx [0.42, 0.29, 0.29]. \quad (5)$$

By adding an entropy regularizer to the utilities,  $\tau\mathcal{H}(x_i)$ , we induce a soft-BR. Therefore, the homotopy approach has the added benefit of partially alleviating gradient bias for moderate temperatures. Further empirical analysis of bias can be found in Appx. F.1.

## 4 ADIDAS

In the previous section, we laid out the conceptual approach we take and identified bias as a potential issue to scaling up computation with Monte Carlo approximation. Here, we inspect the details of our approach, introduce further modifications to reduce the issue of bias, and present our resulting algorithm ADIDAS. Finally, we discuss the advantages of our approach for scaling to large games.

### 4.1 Deviation Incentive Gradient

Regularizing the utilities with weighted Shannon entropy,  $u_k^\tau(\mathbf{x}) = u_k(\mathbf{x}) + S_k^\tau(x_k, x_{-k})$ , where  $S_k^\tau(x_k, x_{-k}) = -\tau \sum_{a_k} x_{ka_k} \ln(x_{ka_k})$ , leads to the following average deviation incentive gradient derived in Appx. E where  $\text{BR}_j = \text{softmax}(\nabla_{x_j}^j / \tau)$  and  $\text{diag}(v)$  creates a diagonal matrix with  $v$  on the diagonal:

$$\nabla_{x_i} \mathcal{L}_{adi}^\tau(\mathbf{x}) = - \overbrace{(\nabla_{x_i}^i - \tau(\ln(x_i) + 1))}^{\text{policy gradient}} + \sum_{j \neq i} \left[ J_{x_i}(\text{BR}_j)^\top (\nabla_{x_j}^j - \tau(\ln(\text{BR}_j) + 1)) + H_{ij}^j(\text{BR}_j - x_j) \right] \quad (6)$$

$$\text{with } J_{x_i}(\text{BR}_j) = \frac{1}{\tau} (\text{diag}(\text{BR}_j) - \text{BR}_j \text{BR}_j^\top) H_{ji}^j. \quad (7)$$

In the limit,  $\nabla_{x_i} \mathcal{L}_{adi}^\tau(\mathbf{x}) \xrightarrow{\tau \rightarrow 0^+} -\nabla_{x_i}^i + \sum_{j \neq i} H_{ij}^j(\text{BR}_j - x_j)$ . The first term is recognized as player  $i$ ’s payoff or *policy* gradient. The second term is a correction that accounts for the other players’

incentives to exploit player  $i$  through a strategy deviation. Each  $H_{ij}^j$  approximates player  $j$ 's payoffs in the bimatrix game between players  $i$  and  $j$ . Recall from the preliminaries that in a polymatrix game, these matrices capture the game exactly. We also explore an adaptive Tsallis entropy in Appx. E.

## 4.2 Amortized Estimates with Historical Play

Section 3.2 discusses the bias that can be introduced when best responding to sampled joint play and how the annealing process of the homotopy method helps alleviate it by softening the BR operator with entropy regularization. To reduce the bias further, we could evaluate more samples from  $\mathbf{x}$ , however, this increases the required computation. Alternatively, assuming strategies have changed minimally over the last few updates (i.e.,  $\mathbf{x}^{(t-2)} \approx \mathbf{x}^{(t-1)} \approx \mathbf{x}^{(t)}$ ), we can instead reuse historical play to improve estimates. We accomplish this by introducing an auxiliary variable  $y_i$  that computes an exponentially averaged estimate of each player  $i$ 's payoff gradient  $\nabla_{x_i}^i$  throughout the descent similarly to Sutton et al. [56]. We also use  $y_i$  to compute an estimate of ADI,  $\hat{\mathcal{L}}_{adi}^\tau$ , as follows:

$$\hat{\mathcal{L}}_{adi}^\tau(\mathbf{x}, \mathbf{y}) = \sum_k y_k^\top (\hat{\text{BR}}_k - x_k) + S_k^\tau(\hat{\text{BR}}_k, x_{-k}) - S_k^\tau(x_k, x_{-k}) \quad (8)$$

where  $\hat{\text{BR}}_k = \arg \max_{z_k \in \Delta^{m_k-1}} y_k^\top z_k + S_k^\tau(z_k, x_{-k})$  is computed with  $y_k$  instead of  $\nabla_{x_k}^k$ . Likewise, replace all  $\nabla_{x_k}^k$  with  $y_k$  and  $\text{BR}_k$  with  $\hat{\text{BR}}_k$  in equations (6) and (7) when computing the gradient:

## 4.3 Putting It All Together

### Algorithm 2 ADIDAS

---

```

1: Given: Strategy learning rate  $\eta_x$ , auxiliary learning rate  $\eta_y$ ,
   initial temperature  $\tau$  ( $= 100$ ), ADI threshold  $\epsilon$ , total iterations  $T$ ,
   simulator  $\mathcal{G}_i$  that returns player  $i$ 's payoff given a joint action.
2:  $\mathbf{x} \leftarrow \{\frac{1}{m_i} \mathbf{1} \forall i\}$ 
3:  $\mathbf{y} \leftarrow \{\mathbf{0} \forall i\}$ 
4: for  $t = 1 : T$  do
5:    $a_i \sim x_i \forall i$ 
6:   for  $i \in \{1, \dots, n\}$  do
7:     for  $j \neq i \in \{1, \dots, n\}$  do
8:        $H_{ij}^i[r, c] \leftarrow \mathcal{G}_i(r, c, a_{-ij}) \forall r \in \mathcal{A}_i, c \in \mathcal{A}_j$ 
9:     end for
10:    end for
11:     $\nabla_{x_i}^i = H_{ij}^i x_j$  for any  $x_j$  (or average the result over all  $j$ )
12:     $y_i \leftarrow y_i - \max(\frac{1}{t}, \eta_y)(\nabla_{x_i}^i - y_i)$ 
13:     $x_i \leftarrow x_i - \eta_x \nabla_{x_i}^i \hat{\mathcal{L}}_{adi}^\tau(\mathbf{x}, \mathbf{y})$  (def. in §4.2 and code in Appx. L)
14:    if  $\hat{\mathcal{L}}_{adi}^\tau(\mathbf{x}, \mathbf{y}) < \epsilon$  (def. in equation (8)) then
15:       $\tau \leftarrow \frac{\tau}{2}$ 
16:    end if
17:  end for
18: return  $\mathbf{x}$ 

```

---

Algorithm 2, ADIDAS, is our final algorithm. ADIDAS attempts to approximate the unique continuum of quantal response equilibria by way of a quasi-stationary process—see Figure 2. Whenever the

algorithm finds a joint strategy  $\mathbf{x}$  exhibiting  $\hat{\mathcal{L}}_{adi}^\tau$  below a threshold  $\epsilon$  for the game regularized with temperature  $\tau$ , the temperature is exponentially reduced (line 15 of ADIDAS) as suggested in [57]. Incorporating stochastic optimization into the process enables scaling the classical homotopy approach to extremely large games (large payoff tensors). At the same time, the homotopy approach selects a unique limiting equilibrium and, symbiotically, helps alleviate gradient bias, further amortized by the reuse of historical play.

*Limitations:* As mentioned earlier, gradient bias precludes a rigorous convergence proof of ADIDAS. However, recent work showed that gradient estimators that are biased, but consistent worked well empirically [18] and follow-up analysis suggests consistency may be an important property [17]. Bias is also being explored in the more complex Riemannian optimization setting where it has been proven that the amount of bias in the gradient shifts the stationary point by a proportional amount [24]. Note that ADIDAS gradients are also consistent in the limit of infinite samples of joint play, and we also find that biased stochastic gradient descent maintains an adequate level of performance for the purpose of our experiments.

No-regret algorithms scale, but have been proven not to converge to Nash [30] and classical solvers [44] converge to Nash, but do not scale. ADIDAS suffers from gradient bias, an issue that may be further mitigated by future research. In this sense, ADIDAS is one of the few, if only, algorithms that can practically approximate Nash in many-player, many-action normal-form games.

Alg Family	Classical	No-Regret	This Work
Convergence to Nash	Yes	No	Yes <sup>†</sup>
Payoffs Queried	$nm^n$	$Tnm^{\ddagger}$	$T(nm)^2$

**Table 2: Comparison of solvers.** <sup>†</sup>See *Limitations* in §4.3 and **Appx. D.2**. <sup>‡</sup>Reduce to  $T$  at the expense of higher variance.

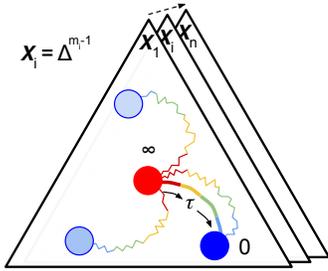
## 4.4 Complexity and Savings

A normal form game may also be represented with a tensor  $U$  in which each entry  $U[i, a_1, \dots, a_n]$  specifies the payoff for player  $i$  under the joint action  $(a_1, \dots, a_n)$ . In order to demonstrate the computational savings of our approach, we evaluate the ratio of the number of entries in  $U$  to the number of entries queried (in the sense of [5, 27, 28]) for computing a single gradient,  $\nabla \mathcal{L}_{adi}^\tau$ . This ratio represents the number of steps that a gradient method can take before it is possible to compute  $\mathcal{L}_{adi}^\tau$  exactly in expectation.

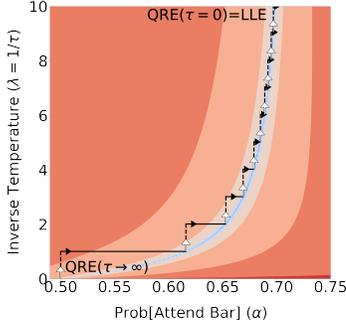
Without further assumptions on the game, the number of entries in a general payoff tensor is  $nm^n$ . In contrast, computing the stochastic deviation incentive gradient requires computing  $H_{ij}^j$  for all  $i, j$  requiring less than  $(nm)^2$  entries<sup>5</sup>. The resulting ratio is  $\frac{1}{n} m^{n-2}$ . For a 7-player, 21-action game, this implies at least 580,000 descent updates can be used by stochastic gradient descent.

If the game is symmetric and we desire a symmetric Nash, the payoff tensor can be represented more concisely with  $\frac{(m+n-1)!}{n!(m-1)!}$  entries (number of multisets of cardinality  $n$  with elements taken from a finite set of cardinality  $m$ ). The number of entries required

<sup>5</sup>Recall  $\nabla_{x_i}^i$  can be computed with  $\nabla_{x_i}^i = H_{ij}^i x_j$  for any  $x_j$ .



(a) ADIDAS pathologies



(b) 10-player, 2-action El Farol homotopy

**Figure 2:** (a) In the presence of multiple equilibria, ADIDAS may fail to follow the path to the uniquely defined Nash due to gradient noise, gradient bias, and a coarse annealing schedule. If these issues are severe, they can cause the algorithm to get stuck at a local optimum of  $\mathcal{L}_{adi}^\tau$ —see Figure 5b in §5.2. (b) Such concerns are minimal for the El Farol Bar stage game by Arthur [4]. The solid black curves represent (biased) descent trajectories while the dashed segments indicate the temperature is being annealed.

for a stochastic gradient is less than  $m^2$ . Again, for a 7-player 21-action game, this implies at least 2,000 update steps. Although there are fewer unique entries in a symmetric game, we are not aware of libraries that allow sparse storage of or efficient arithmetic on such permutation-invariant tensors. ADIDAS can exploit this symmetry.

## 5 EXPERIMENTS

We test the performance of ADIDAS empirically on very large games. We begin by considering Colonel Blotto, a deceptively complex challenge domain still under intense research [9, 13], implemented in OpenSpiel [39]. For reference, both the 3 and 4-player variants we consider are an order of magnitude ( $> 20\times$ ) larger than the largest games explored in [53]. We find that no-regret approaches as well as existing methods from Gambit [44] begin to fail at this scale, whereas ADIDAS performs consistently well. At the same time, we empirically validate our design choice regarding amortizing gradient estimates (§4.2). Finally, we end with our most challenging experiment, the approximation of a unique Nash of a 7-player, 21-action ( $>$  billion outcome) Diplomacy meta-game.

We use the following notation to indicate variants of the algorithms compared in Table 3. A  $y$  superscript prefix, e.g.,  ${}^y\text{QRE}$ , indicates the estimates of payoff gradients are amortized using historical play; its absence indicates that a fresh estimate is used instead.  $\bar{x}_t$  indicates that the average deviation incentive reported is for the average of  $x^{(t)}$  over learning. A subscript of  $\infty$  indicates best responses are computed with respect to the true expected payoff gradient (infinite samples). A superscript *auto* indicates the temperature  $\tau$  is annealed according to line 15 of Algorithm 2. An  $s$  in parentheses indicates lines 5-10 of ADIDAS are repeated  $s$  times, and the resulting  $H_{ij}^i$ 's are averaged for a more accurate estimate. Each game is solved on 1 CPU, except Diplomacy (see Appx. A).

FTRL	Simultaneous Gradient Ascent
RM	Regret-Matching [10]
ATE	ADIDAS with Tsallis (Appx. G)
QRE	ADIDAS with Shannon

**Table 3: Algorithms**

$\eta_x$	$10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}$
$\eta_x^{-1} \cdot \eta_y$	1, 10, 100
$\tau$	0.0, 0.01, 0.05, 0.10
$\Pi_\Delta(\nabla \mathcal{L}_{adi})$	Boolean
Bregman- $\psi(x)$	$\{\frac{1}{2}\ x\ ^2, -\mathcal{H}(x)\}$
$\epsilon$	0.01, 0.05

**Table 4: Hyperparameter Sweeps**

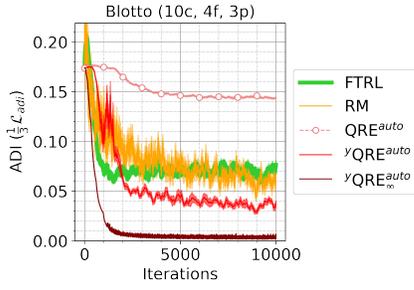
Sweeps are conducted over whether to project gradients onto the simplex ( $\Pi_\Delta(\nabla \mathcal{L}_{adi})$ ), whether to use a Euclidean projection or entropic mirror descent [6] to constrain iterates to the simplex, and over learning rates. Averages over 10 runs of the best hyperparameters are then presented<sup>6</sup> except for Diplomacy for which we present all settings attempted (more in Appx. I.2). Performance is measured by  $\mathcal{L}_{adi}$ , a.k.a. NashConv [40]. For symmetric games, we enforce searching for a symmetric equilibrium (see Appx. C).

For sake of exposition, we do not present all baselines in all plots, however, we include the full suite of comparisons in the appendix. Our experiments demonstrate that without any additional prior information on the game, ADIDAS is the only practical approach for approximating a Nash equilibrium over many-players and many-actions. We argue this by systematically ruling out other approaches on a range of domains. For example, in Figure 3, RM reduces ADI adequately in Blotto. We do not present RM with improvements in Figure 3 such as using exact expectations,  $\text{RM}_\infty$ , or averaging its iterates,  $\text{RM}(\bar{x}_t)$ , because we show that both these fail to save RM on the GAMUT game in Figure 4. In other words, we do not present baselines that are unnecessary for logically supporting the claim above. Code is available at [github.com/deepmind/open\\_spiel](https://github.com/deepmind/open_spiel) [39].

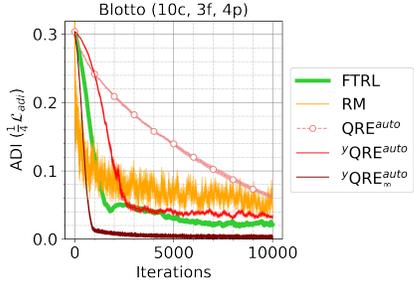
### 5.1 Med-Scale re. §4.4

Govindan-Wilson is considered a state-of-the-art Nash solver, but it does not scale well to large games. For example, on a symmetric, 4-player Blotto game with 66 actions (10 coins, 3 fields), GW, as

<sup>6</sup>Best hyperparameters are used because we expect running ADIDAS with multiple hyperparameter settings in parallel to be a pragmatic approach to approximating Nash.



(a) 3-player, 286-action Blotto



(b) 4-player, 66-action Blotto

**Figure 3: Amortizing estimates of joint play using  $y$  can reduce gradient bias, further improving performance (e.g., compare  $\text{QRE}^{\text{auto}}$  to  $y\text{QRE}^{\text{auto}}$  in (a) or (b)).**

implemented in Gambit, is estimated to take 53,000 hours<sup>7</sup>. Of the solvers implemented in Gambit, none finds a symmetric Nash equilibrium within an hour<sup>8</sup>. Of those, `gambit-logit` [57] is expected to scale most gracefully. Experiments from the original paper are run on maximum 5-player games (2-actions per player) and 20-action games (2-players), so the 4-player, 66-action game is well outside the original design scope. Attempting to run `gambit-logit` anyways with a temperature  $\tau = 1$  returns an approximate Nash with  $\mathcal{L}_{\text{adi}} = 0.066$  after 101 minutes. In contrast, Figure 3b shows ADIDAS achieves a lower ADI in  $\approx 3$  minutes.

*Auxiliary  $y$  re. §4.2.* The introduction of auxiliary variables  $y_i$  are supported by the results in Figure 3— $y\text{QRE}^{\text{auto}}$  significantly improves performance over  $\text{QRE}^{\text{auto}}$  and with low algorithmic cost.

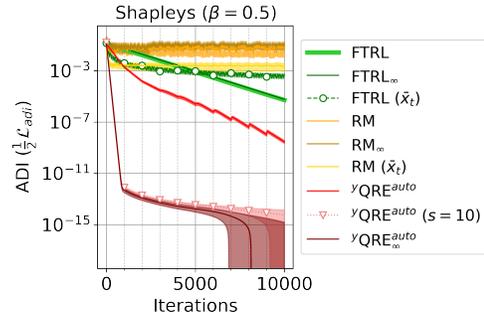
*No-regret, No-convergence re. §4.3.* In Figure 3, FTRL and RM achieve low ADI quickly in some cases. FTRL has recently been proven not to converge to Nash, and this is suggested to be true of no-regret algorithms in general [30, 47]. Before proceeding, we demonstrate empirically in Figure 4 that FTRL and RM fail on games where ADIDAS significantly reduces ADI. Note that GAMUT (D7) was highlighted as a particularly challenging problem for Nash solvers in [53].

## 5.2 Large-Scale

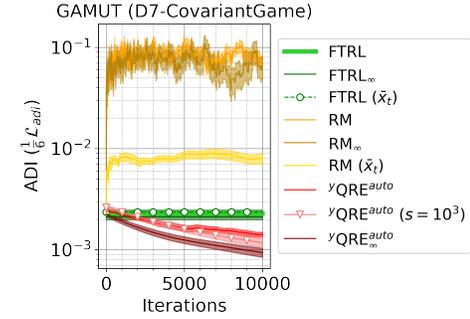
Figure 5 demonstrates an empirical game theoretic analysis [36, 61, 62] of a large symmetric 7-player Diplomacy meta-game where each player elects 1 of 5 trained bots to play on their behalf. Each

<sup>7</sup>Public correspondence with primary `gambit` developer [link].

<sup>8</sup>`gambit-enumpoly` returns several non-symmetric, pure Nash equilibria. Solvers listed in Appx. H.2. Symmetric equilibria are necessary for ranking in symmetric meta-games.



(a) 2-player, 3-action modified Shapley's

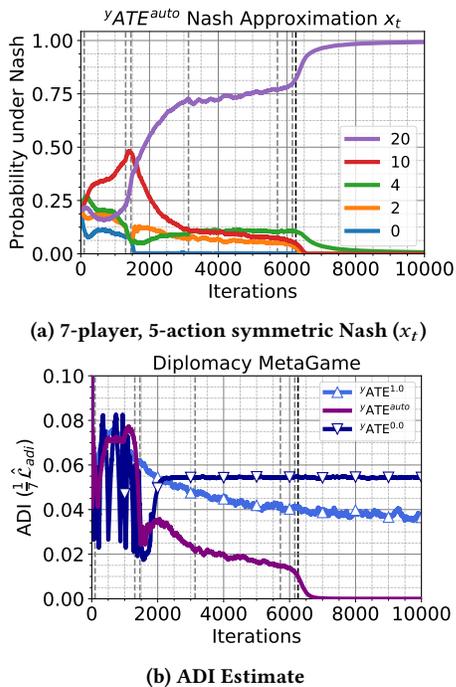


(b) 6-player, 5-action GAMUT-D7

**Figure 4: ADIDAS reduces  $\mathcal{L}_{\text{adi}}$  in both these nonsymmetric games. In contrast, regret matching stalls or diverges in game (a) and diverges in game (b). FTRL makes progress in game (a) but stalls in game (b). In game (a), created by Ostrovski and van Strien [51], to better test performance,  $x$  is initialized randomly rather than with the uniform distribution because the Nash is at uniform.**

bot represents a snapshot taken from an RL training run on Diplomacy [2]. In this case, the expected value of each entry in the payoff tensor represents a winrate. Each entry can only be estimated by simulating game play, and the result of each game is a Bernoulli random variable (ruling out deterministic approaches, e.g., `gambit`). To estimate winrate within 0.01 (ADI within 0.02) of the true estimate with probability 95%, a Chebyshev bound implies more than 223 samples are needed. The symmetric payoff tensor contains 330 unique entries, requiring over 73 thousand games in total. ADIDAS achieves near zero ADI in less than 7 thousand iterations with 50 samples of joint play per iteration ( $\approx 5 \times$  the size of the tensor).

*Continuum of QREs approaching LLE.* The purpose of this work is to approximate a unique Nash (the LLE) which ADIDAS is designed to do, however, the approach ADIDAS takes of attempting to track the continuum of QREs (or the continuum defined by the Tsallis entropy) allows returning these intermediate QRE strategies which may be of interest. Access to these intermediate approximations can be useful when a game playing program cannot wait for ADIDAS's final output to play a strategy, for example, in online play. Interestingly, human play appears to track the continuum of QREs in some cases where the human must both learn about the game (rules, payoffs, etc.) whilst also evolving their strategy [45]. Notice



**Figure 5: (a) Evolution of the symmetric Nash approximation returned by ADIDAS for the 7-player Diplomacy meta-game that considers a subset  $\{0, 2, 4, 10, 20\}$  of the available 21 bots; (b) ADI estimated from auxiliary variable  $y_t$ . Black vertical lines indicate the temperature  $\tau$  was annealed.**

in Figure 5 that the trajectory of the Nash approximation is not monotonic; for example, see the kink around 2000 iterations where bots 10 and 20 swap rank. The continuum of QRE’s from  $\tau = \infty$  to  $\tau = 0$  is known to be complex providing further reason to carefully estimate ADI and its gradients.

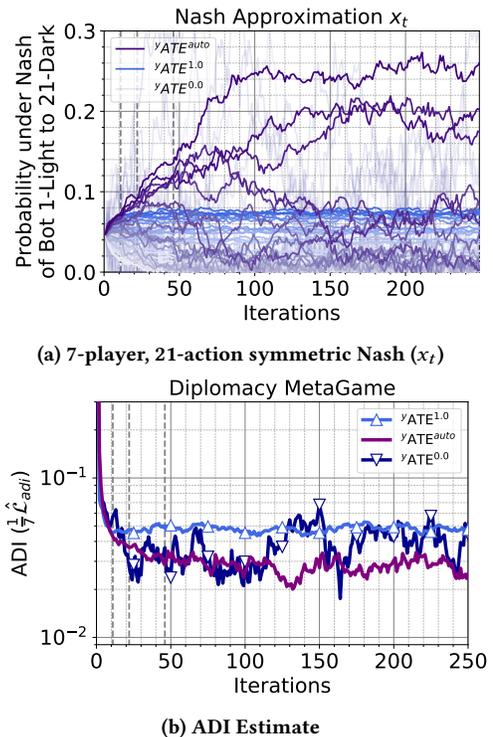
*Convergence to a Local Optimum.* One can also see from Figure 5b that  $y_{ATE^{0.0}}$  has converged to a suboptimal local minimum in the energy landscape. This is likely due to the instability and bias in the gradients computed without any entropy bonus; notice the erratic behavior of its ADI within the first 2000 iterations.

### 5.3 Very Large-Scale re. §4.4

Finally, we repeat the above analysis with all 21 bots. To estimate winrate within 0.015 (ADI within 0.03) of the true estimate with probability 95%, a Chebyshev bound implies approximately 150 samples are needed. The symmetric payoff tensor contains 888,030 unique entries, requiring over 100 million games in total. Note that ignoring the symmetry would require simulating  $150 \times 21^7 \approx 270$  billion games and computing over a trillion payoffs ( $\times 7$  players). Simulating all games, as we show, is unnecessarily wasteful, and just storing the entire payoff tensor in memory, let alone computing with it would be prohibitive without special permutation-invariant data structures ( $\approx 50$ GB with float32). In Figure 6a, ADIDAS with  $\eta_x = \eta_y = 0.1$  and  $\epsilon = 0.001$  achieves a stable ADI below 0.03 in less than 100 iterations with 10 samples of joint play per iteration and each game repeated 7 times ( $< 2.5\%$  of the games run by the naive

alternative). As expected, bots later in training (darker lines) have higher mass under the Nash distribution computed by  $y_{ATE^{auto}}$ . Runtime is discussed in Appx. A.

*Importance of Entropy Bonus.* Figure 6a shows how the automated annealing mechanism ( $y_{ATE^{auto}}$ ) seeks to maintain entropy regularization near a “sweet spot” —too little entropy ( $y_{ATE^{0.0}}$ ) results in an erratic evolution of the Nash approximation and too much entropy ( $y_{ATE^{1.0}}$ ) prevents significant movement from the initial uniform distribution. Figure 6b shows that ADIDAS with the automated annealing mechanism meant to trace the QRE continuum achieves a lower ADI than its fixed temperature variants.



**Figure 6: (a) Evolution of the symmetric Nash approximation returned by ADIDAS for the 7-player, 21-bot Diplomacy meta-game; (b) ADI estimated from auxiliary variable  $y_t$ . Black vertical lines indicate the temperature  $\tau$  was annealed.**

In the appendix, we perform additional ablation studies (e.g., no entropy, annealing), measure accuracy of  $\hat{L}_{adi}^\tau$ , compare against more algorithms on other domains, and consider Tsallis entropy.

## 6 CONCLUSION

Existing algorithms either converge to Nash, but do not scale to large games or scale to large games, but do not converge to Nash. We proposed an algorithm to fill this void that queries necessary payoffs through sampling, obviating storing the full payoff tensor in memory. ADIDAS is principled and shown empirically to approximate Nash in large-normal form games.

## REFERENCES

- [1] AmirMahdi Ahmadinejad, Sina Dehghani, MohammadTaghi Hajiaghayi, Brendan Lucier, Hamid Mahini, and Saeed Seddighin. 2019. From duels to battlefields: Computing equilibria of Blotto and other games. *Mathematics of Operations Research* 44, 4 (2019), 1304–1325.
- [2] Thomas Anthony, Tom Eccles, Andrea Tacchetti, János Kramár, Ian Gemp, Thomas C Hudson, Nicolas Porcel, Marc Lanctot, Julien Pérolat, Richard Everett, et al. 2020. Learning to Play No-Press Diplomacy with Best Response Policy Iteration. In *Advances in Neural Information Processing Systems*.
- [3] Ayala Arad and Ariel Rubinstein. 2012. Multi-dimensional iterative reasoning in action: The case of the Colonel Blotto game. *Journal of Economic Behavior & Organization* 84, 2 (2012), 571–585.
- [4] W Brian Arthur. 1994. Complexity in economic theory: Inductive reasoning and bounded rationality. *The American Economic Review* 84, 2 (1994), 406–411.
- [5] Yakov Babichenko. 2016. Query complexity of approximate Nash equilibria. *Journal of the ACM (JACM)* 63, 4 (2016), 36:1–36:24.
- [6] Amir Beck and Marc Teboulle. 2003. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters* 31, 3 (2003), 167–175.
- [7] Soheil Behnezhad, Avrim Blum, Mahsa Derakhshan, MohammadTaghi Hajiaghayi, Mohammad Mahdian, Christos H Papadimitriou, Ronald L Rivest, Saeed Seddighin, and Philip B Stark. 2018. From battlefields to elections: Winning strategies of Blotto and auditing games. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2291–2310.
- [8] Soheil Behnezhad, Avrim Blum, Mahsa Derakhshan, MohammadTaghi Hajiaghayi, Christos H Papadimitriou, and Saeed Seddighin. 2019. Optimal strategies of Blotto games: Beyond convexity. In *Proceedings of the 2019 ACM Conference on Economics and Computation*. 597–616.
- [9] Soheil Behnezhad, Sina Dehghani, Mahsa Derakhshan, MohammadTaghi Hajiaghayi, and Saeed Seddighin. 2017. Faster and simpler algorithm for optimal strategies of Blotto game. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 369–375.
- [10] David Blackwell et al. 1956. An analog of the minimax theorem for vector payoffs. *Pacific J. Math.* 6, 1 (1956), 1–8.
- [11] Avrim Blum and Yishay Mansour. 2007. *Learning, Regret Minimization, and Equilibria*. Cambridge University Press, 79–102. <https://doi.org/10.1017/CBO9780511800481.006>
- [12] Ben Blum, Christian R Shelton, and Daphne Koller. 2006. A continuation method for Nash equilibria in structured games. *Journal of Artificial Intelligence Research* 25 (2006), 457–502.
- [13] Enric Boix-Adserà, Benjamin L Edelman, and Siddhartha Jayanti. 2020. The Multiplayer Colonel Blotto Game. In *Proceedings of the 21st ACM Conference on Economics and Computation*. 47–48.
- [14] Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. 2015. Heads-up Limit Hold'em Poker is Solved. *Science* 347, 6218 (January 2015), 145–149.
- [15] Noam Brown, Anton Bakhtin, Adam Lerer, and Qucheng Gong. 2020. Combining Deep Reinforcement Learning and Search for Imperfect-Information Games. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 17057–17069. <https://proceedings.neurips.cc/paper/2020/file/c61f571dbd2fb949d3fe5ae1608dd48b-Paper.pdf>
- [16] Noam Brown and Tuomas Sandholm. 2017. Superhuman AI for Heads-up No-limit Poker: Libratus beats top professionals. *Science* 360, 6385 (December 2017).
- [17] Jie Chen and Ronny Luss. 2018. Stochastic gradient descent with biased but consistent gradient estimators. *arXiv preprint arXiv:1807.11880* (2018).
- [18] Jie Chen, Tengfei Ma, and Cao Xiao. 2018. Fastgcn: fast learning with graph convolutional networks via importance sampling. *arXiv preprint arXiv:1801.10247* (2018).
- [19] Xi Chen and Xiaotie Deng. 2006. Settling the complexity of two-player Nash equilibrium. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*. IEEE, 261–272.
- [20] Xi Chen, Xiaotie Deng, and Shang-Hua Teng. 2009. Settling the complexity of computing two-player Nash equilibria. *Journal of the ACM (JACM)* 56, 3 (2009), 1–57.
- [21] Constantinos Daskalakis, Paul W Goldberg, and Christos H Papadimitriou. 2009. The complexity of computing a Nash equilibrium. *SIAM J. Comput.* 39, 1 (2009), 195–259.
- [22] Argyrios Deligkas, John Fearnley, Tobenna Peter Igwe, and Rahul Savani. 2016. An empirical study on computing equilibria in polymatrix games. *arXiv preprint arXiv:1602.06865* (2016).
- [23] Argyrios Deligkas, John Fearnley, Rahul Savani, and Paul Spirakis. 2017. Computing approximate Nash equilibria in polymatrix games. *Algorithmica* 77, 2 (2017), 487–514.
- [24] Alain Durmus, Pablo Jiménez, Éric Moulines, Salem Said, and Hoi-To Wai. 2020. Convergence analysis of Riemannian stochastic approximation schemes. *arXiv preprint arXiv:2005.13284* (2020).
- [25] Kousha Etessami and Mihalis Yannakakis. 2010. On the complexity of Nash equilibria and other fixed points. *SIAM J. Comput.* 39, 6 (2010), 2531–2597.
- [26] Francisco Facchinei and Jong-Shi Pang. 2007. *Finite-dimensional variational inequalities and complementarity problems*. Springer Science & Business Media.
- [27] John Fearnley, Martin Gairing, Paul W Goldberg, and Rahul Savani. 2015. Learning equilibria of games via payoff queries. *The Journal of Machine Learning Research* 16, 1 (2015), 1305–1344.
- [28] John Fearnley and Rahul Savani. 2016. Finding approximate Nash equilibria of bimatrix games via payoff queries. *ACM Transactions on Economics and Computation (TEAC)* 4, 4 (2016), 25:1–25:19.
- [29] Mark Fey. 2012. Symmetric games with only asymmetric equilibria. *Games and Economic Behavior* 75, 1 (2012), 424–427.
- [30] Lampros Flokas, Emmanouil-Vasileios Vlatakis-Gkaragkounis, Thanasis Lianeas, Panayotis Mertikopoulos, and Georgios Piliouras. 2020. No-regret learning and mixed Nash equilibria: They do not mix. *arXiv preprint arXiv:20.09514* (2020).
- [31] Drew Fudenberg, David K Levine, et al. 1998. *The Theory of Learning in Games*. MIT Press Books 1 (1998).
- [32] Paul W Goldberg, Rahul Savani, Troels Bjerre Sørensen, and Carmine Ventre. 2013. On the approximation performance of fictitious play in finite games. *International Journal of Game Theory* 42, 4 (2013), 1059–1083.
- [33] Srihari Govindan and Robert Wilson. 2003. A global Newton method to compute Nash equilibria. *Journal of Economic Theory* 110, 1 (2003), 65–86.
- [34] Srihari Govindan and Robert Wilson. 2004. Computing Nash equilibria by iterated polymatrix approximation. *Journal of Economic Dynamics and Control* 28, 7 (2004), 1229–1241.
- [35] Jonathan Gray, Adam Lerer, Anton Bakhtin, and Noam Brown. 2020. Human-Level Performance in No-Press Diplomacy via Equilibrium Search. *arXiv preprint arXiv:2010.02923* (2020).
- [36] Patrick R Jordan, Christopher Kiekintveld, and Michael P Wellman. 2007. Empirical game-theoretic analysis of the TAC supply chain game. In *Proceedings of the 6th International Joint Conference on Autonomous Agents and Multiagent Systems*. 1–8.
- [37] Anatoli Juditsky, Arkadi Nemirovski, and Claire Tauvel. 2011. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems* 1, 1 (2011), 17–58.
- [38] Galina M Korpelevich. 1976. The extragradient method for finding saddle points and other problems. *Matecon* 12 (1976), 747–756.
- [39] Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay, Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, Daniel Hennes, Dustin Morrill, Paul Muller, Timo Ewalds, Ryan Faulkner, János Kramár, Bart De Vylder, Brennan Saeta, James Bradbury, David Ding, Sebastian Borgeaud, Matthew Lai, Julian Schrittwieser, Thomas Anthony, Edward Hughes, Ivo Danihelka, and Jonah Ryan-Davis. 2019. OpenSpiel: A Framework for Reinforcement Learning in Games. *CoRR abs/1908.09453* (2019). [arXiv:1908.09453 \[cs.LG\]](http://arxiv.org/abs/1908.09453) <http://arxiv.org/abs/1908.09453>
- [40] Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. 2017. A unified game-theoretic approach to multiagent reinforcement learning. In *Advances in Neural Information Processing Systems*. 4190–4203.
- [41] Carlton Lemke and Joseph Howson, Jr. 1964. Equilibrium points of bimatrix games. *Journal of the Society for industrial and Applied Mathematics* 12, 2 (1964), 413–423.
- [42] Edward Lockhart, Marc Lanctot, Julien Pérolat, Jean-Baptiste Lespiau, Dustin Morrill, Finbarr Timbers, and Karl Tuyls. 2019. Computing Approximate Equilibria in Sequential Adversarial Games by Exploitability Descent. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI)*.
- [43] Richard D McKelvey and Andrew McLennan. 1996. Computation of equilibria in finite games. *Handbook of Computational Economics* 1 (1996), 87–142.
- [44] Richard D McKelvey, Andrew M McLennan, and Theodore L Turocy. 2016. *Gambit: Software tools for game theory*, version 16.0.1.
- [45] Richard D McKelvey and Thomas R Palfrey. 1995. Quantal response equilibria for normal form games. *Games and Economic Behavior* 10, 1 (1995), 6–38.
- [46] H Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. 2003. Planning in the presence of cost functions controlled by an adversary. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*. 536–543.
- [47] Panayotis Mertikopoulos, Christos Papadimitriou, and Georgios Piliouras. 2018. Cycles in adversarial regularized learning. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2703–2717.
- [48] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. 2017. The Numerics of GANs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (Long Beach, California, USA) (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 1823–1833.
- [49] Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. 2017. DeepStack: Expert-level artificial intelligence in Heads-up No-limit Poker. *Science* 358, 6362 (October 2017).
- [50] Eugene Nudelman, Jennifer Wortman, Yoav Shoham, and Kevin Leyton-Brown. 2004. Run the Gamut: A comprehensive approach to evaluating game-theoretic algorithms. In *AAMAS*, Vol. 4. 880–887.

- [51] Georg Ostrovski and Sebastian van Strien. 2013. Payoff performance of fictitious play. *arXiv preprint arXiv:1308.4049* (2013).
- [52] Julien Perolat, Remi Munos, Jean-Baptiste Lespiau, Shayegan Omidshafiei, Mark Rowland, Pedro Ortega, Neil Burch, Thomas Anthony, David Balduzzi, Bart De Vylder, et al. 2020. From Poincaré Recurrence to Convergence in Imperfect Information Games: Finding Equilibrium via Regularization. *arXiv preprint arXiv:2002.08456* (2020).
- [53] Ryan Porter, Eugene Nudelman, and Yoav Shoham. 2008. Simple search methods for finding a Nash equilibrium. *Games and Economic Behavior* 63, 2 (2008), 642–662.
- [54] Tuomas Sandholm, Andrew Gilpin, and Vincent Conitzer. 2005. Mixed-integer programming methods for finding Nash equilibria. In *AAAI*. 495–501.
- [55] Yoav Shoham and Kevin Leyton-Brown. 2009. *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge University Press.
- [56] Richard S Sutton, Csaba Szepesvári, and Hamid Reza Maei. 2008. A convergent  $O(n)$  algorithm for off-policy temporal-difference learning with linear function approximation. *Advances in Neural Information Processing Systems* 21, 21 (2008), 1609–1616.
- [57] Theodore L Turocy. 2005. A dynamic homotopy interpretation of the logistic quantal response equilibrium correspondence. *Games and Economic Behavior* 51, 2 (2005), 243–263.
- [58] Gerard van der Laan, AJJ Talman, and L Van der Heyden. 1987. Simplicial variable dimension algorithms for solving the nonlinear complementarity problem on a product of unit simplices using a general labelling. *Mathematics of Operations Research* 12, 3 (1987), 377–397.
- [59] Oriol Vinyals, Igor Babuschkin, Junyoung Chung, Michael Mathieu, Max Jaderberg, Wojciech M Czarnecki, Andrew Dudzik, Aja Huang, Petko Georgiev, Richard Powell, et al. 2019. Alphastar: Mastering the real-time strategy game starcraft ii. *DeepMind blog* (2019), 2.
- [60] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575, 7782 (2019), 350–354.
- [61] Elaine Wah, Sébastien Lahaie, and David M Pennock. 2016. An Empirical Game-Theoretic Analysis of Price Discovery in Prediction Markets.. In *IJCAI*. 510–516.
- [62] Michael P Wellman. 2006. Methods for empirical game-theoretic analysis. In *AAAI*. 1552–1556.
- [63] Duncan Whitehead. 2008. The El Farol bar problem revisited: Reinforcement learning in a potential game. *ESE Discussion Papers* 186 (2008).

# Appendices

## CONTENTS

A	Runtime	12
B	Two vs More Than Two Player Games	12
C	Symmetric Nash for Symmetric Games	12
D	Convergence of ADIDAS	12
D.1	Convergence Warm-up: Full Access to In-Memory Payoff Tensor	12
D.2	Convergence Sketch: Sampling the Payoff Tensor	13
E	Deviation Incentive Gradient	13
E.1	Gradient of Utility	13
E.2	Tsallis-Entropy	14
E.3	Shannon Entropy	17
F	Ablations	18
F.1	Bias re. §3.2+§4.1	18
F.2	Auxiliary $y$ re. §4.2	18
F.3	Annealing $\tau$ re. §3.1	18
F.4	Convergence re. §4.3	18
F.5	ADI stochastic estimate	19
G	Experiments Repeated with ATE	19
H	Comparison Against Additional Algorithms	22
H.1	ED and FP Fail	22
H.2	Gambit Solvers	23
I	Additional Game Domains	24
I.1	Diplomacy Experiments - Subsampled Games	24
I.2	Diplomacy Experiments - Empirical Game Theoretic Analysis	25
I.3	El Farol Bar Stage Game	25
J	Description of Domains	25
J.1	Modified Shapley's	25
J.2	Colonel Blotto	25
K	Connections to Other Algorithms	26
K.1	Consensus Algorithm	26
K.2	Exploitability Descent as Extragradient	26
L	Python Code	27

## A RUNTIME

We briefly discussed runtime of ADIDAS in the main body within the context of the Colonel Blotto game. The focus of the paper is on the divide between algorithms that can solve for Nash in any reasonable amount of time (e.g.,  $\approx 3$  minutes) and those that cannot (e.g., GW with 53,000 hours). The modified Shapley’s game and D7-Covariant game from GAMUT are both smaller than the Blotto game, so we omitted a runtime discussion for these.

The Diplomacy experiment required simulating Diplomacy games on a large shared compute cluster with simulated games taking anywhere from 3 minutes to 3 hours. Games were simulated at each iteration of ADIDAS asynchronously using a pool of 1000 workers (4 CPUs per worker, 1 worker per game); the Nash approximate  $x_t$  was updated separately on a single CPU. The main computational bottleneck in this experiment was simulating the games themselves, rather than computing gradients from those games. Therefore, the number of games simulated (entries accessed in the payoff tensor) is a realistic metric of algorithmic efficiency.

## B TWO VS MORE THAN TWO PLAYER GAMES

An  $n$ -player game for all  $n \geq 3$  can be reduced in polynomial time to a 2-player game such that the Nash equilibria of the 2-player game can be efficiently used to compute approximate Nash equilibria of the  $n$ -player game [19, 21, 25].

## C SYMMETRIC NASH FOR SYMMETRIC GAMES

Note that a symmetric Nash equilibrium is guaranteed to exist for a finite, normal-form game [29].

One of the reasons we enforce symmetry is that we had Nash-ranking in mind when designing the algorithm and experiments. In that case, for a symmetric meta-game, we desire a symmetric equilibrium so we have a single ranking to go by for evaluation. If each player, in for example the 7-player Diplomacy meta-game, returned a different distribution at Nash, then we’d have to figure out which player’s side of the Nash to use for ranking.

## D CONVERGENCE OF ADIDAS

We first establish convergence of the simplified algorithm as described in the warm-up and then discuss convergence of the our more sophisticated, scalable algorithm ADIDAS.

### D.1 Convergence Warm-up: Full Access to In-Memory Payoff Tensor

The proof of this simple warm-up algorithm relies heavily on the detailed examination of the continuum of QREs proposed in [45] and further analyzed in [57]. The Theorem presented below is essentially a succinct repetition of one of their results.

**ASSUMPTION 1 (NO PRINCIPAL BRANCHING).** *The continuum of QREs from the uniform Nash to the limiting logit equilibrium is unique and contains no branching points.*

**ASSUMPTION 2 (NO TURNING POINTS).** *The continuum of QREs from the uniform Nash to the limiting logit equilibrium proceeds along a path with monotonically decreasing (increasing)  $\tau(\lambda)$ .*

**ASSUMPTION 3 (BOUNDED SENSITIVITY OF QRE TO TEMPERATURE).** *The shift in location of the QRE is upper bounded by an amount proportional to the increase in inverse temperature:  $\|x_{\lambda+\Delta\lambda}^* - x_\lambda^*\| \leq \sigma\Delta\lambda$ .*

**ASSUMPTION 4 (BOUND ON BOA’S OF QRE’S).** *Under gradient descent dynamics, the basin of attraction for any quantal response equilibrium,  $x_\lambda^* = \text{QRE}_{\tau=\lambda^{-1}}$ , contains a ball of radius  $r$ . Formally, assuming  $x_{t+1} \leftarrow x_t - \eta_t g_t$  with  $g_t = \nabla_x \mathcal{L}_{adi}^\tau(x_t)$ ,  $\eta_t$  a square-summable, not summable step size (e.g.,  $\propto t^{-1}$ ), and given  $x_0 \in B(x_\lambda^*, r)$ , there exists a  $T$  such that  $x_{t \geq T} \in B(x_\lambda^*, \epsilon)$  for any  $\epsilon$ .*

**THEOREM 1.** *Assume the QREs along the homotopy path have bounded sensitivity to  $\lambda$  given by a parameter  $\sigma$  (Assumption 3), and basins of attraction with radii lower bounded by  $r$  (Assumption 4). Let the step size  $\Delta\lambda \leq \sigma(r - \epsilon)$  with tolerance  $\epsilon$ . And let  $T^*$  be the supremum over all  $T$  such that Assumption 4 is satisfied for any inverse temperature  $\lambda \geq \Delta\lambda$ . Then, assuming gradient descent for OPT, Algorithm 1 converges to the limiting logit equilibrium  $x_{\lambda=\infty}^* = x_{\tau=0}^*$  in the limit as  $T_\lambda \rightarrow \infty$ .*

**PROOF.** Recall McKelvey and Palfrey [45] proved there exists a unique continuum of QREs tracing from infinite temperature ( $\lambda = 0$ ) to zero temperature ( $\lambda = \infty$ ) for almost all games. Assumption 4 effectively assumes the game in question is one from that class. Algorithm 1 initializes  $\lambda = 0$  and  $x$  to the uniform distribution which is the exact QRE for that temperature. Next, in step 5, the temperature is annealed by an amount that, by Lemma 3, ensures  $\|x_{\lambda+\Delta\lambda}^* - x\| = \|x_{\lambda+\Delta\lambda}^* - x_\lambda^*\| \leq r - \epsilon$ , where  $r$  is a minimal radius of the basin of attraction for any QRE. Then, in step 6, OPT returns an  $\epsilon$ -approximation,  $x$ , to the new QRE after  $T^*$  steps, which implies  $\|x - x_{\lambda+\Delta\lambda}^*\| \leq \epsilon$ . The proof then continues by induction. The inverse temperature is increased by an amount ensuring then next QRE is within  $r - \epsilon$  of the previous. The current approximation,  $x$  is within  $\epsilon$  of the previous, therefore, it is within  $r - \epsilon + \epsilon = r$  of the next QRE, i.e., it is in its basin of attraction. The inverse temperature  $\lambda$  is always increased by an amount such that the current approximation is always within the boundary of attraction for the next QRE. Therefore, in the limit of infinite annealing steps,  $x$  converges to the QRE with zero temperature, known as the limiting logit equilibrium.  $\square$

## D.2 Convergence Sketch: Sampling the Payoff Tensor

We do not rigorously prove any theoretical convergence result for the stochastic setting. A convergence proof is complicated by the fact that despite our efforts to reduce gradient bias, some bias will always remain. Although we make assumptions that ensure each iterate begins in the basin of attraction of the QRE of interest, even proving convergence of a hypothetically unbiased stochastic gradient descent to that specific local minimum could only be guaranteed with high probability (dependent on step size). Our goal was to outline a sensible argument that ADIDAS would converge to Nash asymptotically. Our claim of convergence stands on the shoulders of the work of McKelvey and Palfrey [45] who proved that there exists a unique path  $P$  of Quantal Response Equilibria (QREs) parameterized by temperature  $\tau$  which begins at the uniform distribution Nash ( $\tau = \infty$ ) and ends at the limiting logit equilibrium ( $\tau = 0$ ). Turocy [57] solves for this path explicitly by solving the associated initial value problem (differential equation) where  $t = \frac{1}{\tau}$  takes the place of the typical independent variable time. By numerically integrating this differential equation with infinitesimally small steps  $dt$ , Turocy [57] can ensure the iterates progress along the path towards the limiting logit equilibrium (LLE). ADIDAS takes a conceptually similar approach. First, it initializes to the uniform equilibrium. Then it takes a small step  $\Delta t$ . In practice, the initial step we take increases  $t$  from 0 to 1, which worked well enough, but one can imagine taking a smaller step, e.g., 0 to  $10^{-9}$ . After such a small step, the QRE of the game with lower temperature will not have moved far from the initial uniform equilibrium. Therefore, we can minimize ADI to solve for the new QRE, thereby recovering to a point on the unique path  $P$ . The fact that we can only access the payoff tensor by samples means that we may need to sample many times ( $s$  times) to obtain an accurate Monte Carlo estimate of the gradient of ADI. By repeating this process of decaying the temperature ( $\tau_k > \tau_{k+1} \Leftrightarrow t_k < t_{k+1}$ ) and recovering the new QRE with gradient descent (possibly  $n_k$  steps) on ADI ( $x_t = x(\tau_k) \rightarrow x_{t+n_k} = x(\tau_{k+1})$ ), we too can follow  $P$ . In the limit as  $s$ ,  $n_k$ , and  $N = \sum_k n_k$  go to infinity and  $\Delta t$  goes to zero, the issues identified in Figure 2a are mitigated and we recover the LLE. Note,  $n_k$  is effectively increased by reducing  $\epsilon$  in Algorithm 2. We claim "ADIDAS is the first that can approximate Nash in large many-player, many-action normal-form games" because, in principle, it is technically sound according to the argument just presented but also efficient (does not require infinite samples in practice) as demonstrated empirically in our experiments. Note that because we only argue ADIDAS is asymptotically convergent (we provide no convergence rates), we do not contradict any Nash complexity results.

## E DEVIATION INCENTIVE GRADIENT

We now provide the general form for the ADI gradient for normal form games.

$$\begin{aligned} \nabla_{x_i} \mathcal{L}_{adi}(\mathbf{x}) &= \nabla_{x_i} [u_i^\tau(\text{BR}(x_{-i}), x_{-i}) - u_i^\tau(x_i, x_{-i})] \\ &\quad + \sum_{j \neq i} \nabla_{x_i} [u_j^\tau(\text{BR}(x_{-j}), x_{-j}) - u_j^\tau(x_j, x_{-j})]. \end{aligned} \quad (9)$$

$$\nabla_{x_i} [u_i^\tau(\text{BR}(x_{-i}), x_{-i}) - u_i^\tau(x_i, x_{-i})] = \underbrace{J_{x_i}(\text{BR}(x_{-i}))^\top}_{\mathbf{0}} (\nabla_{z_i} u_i^\tau(z_i, x_{-i})|_{\text{BR}_i, x_{-i}}) \quad (10)$$

$$\begin{aligned} &+ \sum_{k \neq i} \underbrace{J_{x_i}(x_k)^\top}_{\mathbf{0}} (\nabla_{z_k} u_i^\tau(\text{BR}(x_i), z_{-i})|_{\text{BR}_i, x_{-i}}) \\ &- \nabla_{x_i} u_i^\tau(x_i, x_{-i}) = -\nabla_{x_i} u_i^\tau(x_i, x_{-i}). \end{aligned} \quad (11)$$

$$\nabla_{x_i} [u_j^\tau(\text{BR}(x_{-j}), x_{-j}) - u_j^\tau(x_j, x_{-j})] = J_{x_i}(\text{BR}(x_{-j}))^\top (\nabla_{z_j} u_j^\tau(z_j, x_{-j})|_{\text{BR}_j, x_{-j}}) \quad (12)$$

$$\begin{aligned} &+ \sum_{k \neq j} J_{x_i}(x_k)^\top (\nabla_{z_k} u_j^\tau(\text{BR}(x_{-j}), z_{-j})|_{\text{BR}_j, x_{-j}}) \\ &- \nabla_{x_i} u_j^\tau(x_j, x_{-j}) \\ &= J_{x_i}(\text{BR}(x_{-j}))^\top (\nabla_{z_j} u_j^\tau(z_j, x_{-j})|_{\text{BR}_j, x_{-j}}) \end{aligned} \quad (13)$$

$$+ (\nabla_{z_j} u_j^\tau(\text{BR}(x_{-j}), z_{-j})|_{\text{BR}_j, x_{-j}}) \quad (14)$$

$$- \nabla_{x_i} u_j^\tau(x_j, x_{-j}). \quad (15)$$

For entropy regularized utilities  $u_i^\tau = u_i + S_i^\tau$ , the policy gradient decomposes as

$$\nabla_{x_i} u_j^\tau(x_j, x_{-j}) = \nabla_{x_i} u_j(x_j, x_{-j}) + \nabla_{x_i} S_j^\tau(x_j, x_{-j}). \quad (16)$$

### E.1 Gradient of Utility

Before deriving the ADI gradient, we first show that the partial derivative of the utility can be defined in terms of expected utilities or payoffs. This eases presentation later. For example, in a 3-player game, player 1's expected utility is defined as

$$u_1(a_1, a_2, a_3) = \sum_{a_1 \in \mathcal{A}_1} \sum_{a_2 \in \mathcal{A}_2} \sum_{a_3 \in \mathcal{A}_3} u_1(a_1, a_2, a_3) x_{1a_1} x_{2a_2} x_{3a_3}. \quad (17)$$

Taking the derivative with respect to player 1's strategy  $x_{1a_1}$ , i.e., the probability of player 1 specifically playing action  $a_1$ , we find

$$\frac{\partial u_1}{\partial x_{1a_1}} = \sum_{a_2 \in \mathcal{A}_2} \sum_{a_3 \in \mathcal{A}_3} u_1(a_1, a_2, a_3) x_{2a_2} x_{3a_3} \quad (18)$$

$$= \mathbb{E}_{a_2 \sim x_2, a_3 \sim x_3} [u_1(a_1, a_2, a_3)]. \quad (19)$$

From here, it should be clear that the full vector of partial derivatives, i.e., gradient, can be written as

$$\nabla_{x_1}(u_1) = \mathbb{E}_{a_2 \sim x_2, a_3 \sim x_3} [u_1(a_1, a_2, a_3)] \forall a_1 \in \mathcal{A}_1. \quad (20)$$

The Jacobian can be defined similarly (e.g, consider differentiating this result w.r.t.  $x_{2a_2}$  for all  $a_2$ ).

## E.2 Tsallis-Entropy

First we derive gradients assuming utilities are carefully regularized using a Tsallis entropy bonus,  $S_k^\tau$ , parameterized by *temperature*  $\tau = p \in [0, 1]$ :

$$S_k^\tau(x_k, x_{-k}) = \frac{s_k}{p+1} \left(1 - \sum_m x_{km}^{p+1}\right) = s_k \frac{p}{p+1} \overbrace{\left[ \frac{1}{p} \left(1 - \sum_m x_{km}^{p+1}\right) \right]}^{\text{Tsallis entropy}} \quad (21)$$

where  $s_k = \left(\sum_m (\nabla_{x_{km}}^k)^{1/p}\right)^p = \|\nabla_{x_k}^k\|_{1/p}$ . For Tsallis entropy, we assume payoffs in the game have been offset by a constant so that they are positive.

The coefficients in front of the Tsallis entropy term are chosen carefully such that a *best response* for player  $k$  can be efficiently computed:

$$\text{BR}(x_{-k}) = \arg \max_{z_k \in \Delta} z_k^\top \nabla_{x_k}^k + \frac{s_k}{p+1} \left(1 - \sum_m z_{km}^{p+1}\right). \quad (22)$$

First note that the maximization problem above is strictly concave for  $s_k > 0$  and  $p \in (0, 1]$ . If these assumptions are met, then any maximum is a unique global maximum. This is a constrained optimization problem, so in general the gradient need not be zero at the global optimum, but in this case it is. We will find a critical point by setting the gradient equal to zero and then prove that this point lies in the feasible set (the simplex) and satisfies second order conditions for optimality.

$$\nabla_{x_k} u_k^\tau(x_k, x_{-k}) = \nabla_{x_k}^k - \|\nabla_{x_k}^k\|_{1/p} x_k^p = 0 \quad (23)$$

$$\implies \text{BR}(x_{-k}) = \left[ \frac{\nabla_{x_k}^k}{\|\nabla_{x_k}^k\|_{1/p}} \right]^{\frac{1}{p}} = \left[ \frac{\nabla_{x_k}^k}{s_k} \right]^{\frac{1}{p}} = \frac{(\nabla_{x_k}^k)^{\frac{1}{p}}}{\|\nabla_{x_k}^k\|_{1/p}^{1/p}} = \frac{(\nabla_{x_k}^k)^{\frac{1}{p}}}{\sum_m (\nabla_{x_{km}}^k)^{\frac{1}{p}}} \in \Delta. \quad (24)$$

The critical point is on the simplex as desired. Furthermore, the Hessian at the critical point is negative definite,  $H(\text{BR}) = -ps_k \text{diag}(\text{BR}^{-1}) < 0$ , so this point is a local maximum (and by strict concavity, a unique global maximum).

If the original assumptions are not met and  $s_k = 0$ , then this necessarily implies  $u_k^\tau(x_k, x_{-k}) = 0$  for all  $x_k$ . As all actions achieve equal payoff, we define the best response in this case to be the uniform distribution. Likewise, if  $p = 0$ , then the Tsallis entropy regularization term disappears ( $1 - \sum_m x_{km} = 0$ ) and the best response is the same as for the unregularized setting. Note in the unregularized setting, we define the best response to be a mixed strategy over all actions achieving the maximal possible utility.

*E.2.1 Gradients.* We now derive the necessary derivatives for computing the deviation incentive gradient.

Entropy Gradients.

$$(A) \nabla_{x_i} S_i^\tau(x_i, x_{-i}) = -s_i x_i^p \quad (25)$$

$$\begin{aligned} (B) \nabla_{x_i} S_j &= p \left( \sum_m (\nabla_{x_{jm}}^j)^{1/p} \right)^{p-1} \left( \sum_m \frac{1}{p} (\nabla_{x_{jm}}^j)^{\frac{1}{p}-1} H_{jmi}^j \right) \\ &= \left( \sum_m (\nabla_{x_{jm}}^j)^{1/p} \right)^{p-1} \left( \sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}-1} H_{jmi}^j \right) \\ &= \frac{1}{s_j^{\frac{1}{p}-1}} H_{ij}^j (\nabla_{x_j}^j)^{\frac{1}{p}-1} = H_{ij}^j \text{BR}(x_{-j})^{1-p} \\ &\stackrel{p=1}{=} H_{ij}^j \mathbf{1} \\ &\stackrel{p=\frac{1}{2}}{=} H_{ij}^j \frac{\nabla_{x_j}^j}{s_j} \end{aligned} \quad (26)$$

$$(C) \nabla_{x_i} S_j^\tau(x_j, x_{-j}) = \frac{1}{s_j} S_j^\tau(x_j, x_{-j}) \underbrace{\nabla_{x_i} S_j}_{(B)} \quad (27)$$

Best Response Gradients.

$$\begin{aligned} (D) J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}] &= J_{x_i} [(H_{ij}^j x_i)^{\frac{1}{p}}] \\ &= \frac{1}{p} (\nabla_{x_j}^j)^{\frac{1}{p}-1} \odot H_{ji}^j \end{aligned} \quad (28)$$

where  $\odot$  denotes elementwise multiplication or, more generally, broadcast multiplication. In this case,  $(\nabla_{x_j}^j)^{\frac{1}{p}-1} \in \mathbb{R}^{d_j \times 1}$  is broadcast multiplied by  $H_{ji}^j \in \mathbb{R}^{d_j \times d_i}$  to produce a Jacobian matrix in  $\mathbb{R}^{d_j \times d_i}$ .

$$\begin{aligned} (E) J_{x_i}(\text{BR}(x_{-j})) &= \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}] - [(\nabla_{x_j}^j)^{\frac{1}{p}}] \left[ \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} \right]^2 \nabla_{x_i} \left[ \sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}} \right]^\top \\ &= \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}] - [(\nabla_{x_j}^j)^{\frac{1}{p}}] \left[ \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} \right]^2 \left[ \sum_m J_{x_i} [(\nabla_{x_{jm}}^j)^{\frac{1}{p}}] \right]^\top \\ &= \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}] - [(\nabla_{x_j}^j)^{\frac{1}{p}}] \left[ \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} \right]^2 [\mathbf{1}^\top J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}]] \\ &= \left[ \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} \mathbf{I}_j - [(\nabla_{x_j}^j)^{\frac{1}{p}}] \left[ \frac{1}{\sum_m (\nabla_{x_{jm}}^j)^{\frac{1}{p}}} \right]^2 \mathbf{1}^\top \right] J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}] \\ &= \frac{1}{\|\nabla_{x_j}^j\|_{1/p}^{1/p}} \left[ \mathbf{I}_j - \frac{(\nabla_{x_j}^j)^{\frac{1}{p}}}{\|\nabla_{x_j}^j\|_{1/p}^{1/p}} \mathbf{1}^\top \right] J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}] \\ &= \frac{1}{\|\nabla_{x_j}^j\|_{1/p}^{1/p}} \left[ \mathbf{I}_j - \text{BR}(x_{-j}) \mathbf{1}^\top \right] \underbrace{J_{x_i} [(\nabla_{x_j}^j)^{\frac{1}{p}}]}_{(D)} \\ &= \frac{1}{s_j^{1/p}} \left[ \mathbf{I}_j - \text{BR}(x_{-j}) \mathbf{1}^\top \right] \left[ \frac{1}{p} (\nabla_{x_j}^j)^{\frac{1}{p}-1} \odot H_{ji}^j \right] \end{aligned} \quad (29)$$

*Deviation Incentive Gradient Terms.* Here, we derive each of the terms in the ADI gradient. The numbers left of the equations mark which terms we are computing in section E.

$$\begin{aligned} (11) \nabla_{x_i} [u_i^\tau(\text{BR}(x_{-i}), x_{-i}) - u_i^\tau(x_i, x_{-i})] &= -\nabla_{x_i} u_i^\tau(x_i, x_{-i}) \\ &\stackrel{(16)+(A)}{=} -(\nabla_{x_i}^i - s_i x_i^p). \end{aligned} \quad (30)$$

$$\begin{aligned}
(13) \quad \nabla_{z_j} u_j^\tau(z_j, x_{-j})|_{\text{BR}_j, x_{-j}} &= [\nabla_{z_j} u_j(z_j, x_{-j}) + \underbrace{\nabla_{z_j} S_j^\tau(z_j, x_{-j})}_{(A)}] |_{\text{BR}_j, x_{-j}} \\
&= [\nabla_{z_j}^j - s_j z_j^p] |_{\text{BR}_j, x_{-j}} \\
&= \nabla_{x_j}^j - s_j \text{BR}(x_{-j})^p \\
&= \nabla_{x_j}^j - \nabla_{x_j}^j = 0.
\end{aligned} \tag{31}$$

$$\begin{aligned}
(14) \quad \nabla_{z_i} u_i^\tau(\text{BR}(x_{-j}), z_{-j})|_{\text{BR}_j, x_{-j}} &= [\nabla_{z_i} u_i(\text{BR}(x_{-j}), z_{-j}) + \underbrace{\nabla_{z_i} S_i^\tau(\text{BR}(x_{-j}), z_{-j})}_{(C)}] |_{\text{BR}_j, x_{-j}} \\
&= [H_{ij}^j \text{BR}(x_{-j}) + \frac{1}{s_j} S_j^\tau(\text{BR}(x_{-j}), x_{-j}) H_{ij}^j \text{BR}(x_{-j})^{1-p}] \\
&= H_{ij}^j \text{BR}(x_{-j}) [1 + \frac{1}{s_j} S_j^\tau(\text{BR}(x_{-j}), x_{-j}) \text{BR}(x_{-j})^{-p}]
\end{aligned} \tag{32}$$

$$\begin{aligned}
(15) \quad \nabla_{z_i} u_i^\tau(x_j, z_{-j})|_{x_j, x_{-j}} &= [\nabla_{z_i} u_i(x_j, z_{-j}) + \underbrace{\nabla_{z_i} S_i^\tau(x_j, z_{-j})}_{(C)}] |_{x_j, x_{-j}} \\
&= [H_{ij}^j x_j + \frac{1}{s_j} S_j^\tau(x_j, x_{-j}) H_{ij}^j x_j^{1-p}] \\
&= H_{ij}^j x_j [1 + \frac{1}{s_j} S_j^\tau(x_j, x_{-j}) x_j^{-p}]
\end{aligned} \tag{33}$$

$$\begin{aligned}
(12) \quad \nabla_{x_i} [u_i^\tau(\text{BR}(x_{-j}), x_{-j}) - u_i^\tau(x_j, x_{-j})] & \tag{34} \\
&\stackrel{(13)+(14)-(15)}{=} H_{ij}^j \text{BR}(x_{-j}) [1 + \frac{1}{s_j} S_j^\tau(\text{BR}(x_{-j}), x_{-j}) \text{BR}(x_{-j})^{-p}] - H_{ij}^j x_j [1 + \frac{1}{s_j} S_j^\tau(x_j, x_{-j}) x_j^{-p}] \\
&= H_{ij}^j (\text{BR}(x_{-j}) - x_j) \\
&+ \frac{1}{p+1} (1 - \|\text{BR}(x_{-j})\|_{p+1}^{p+1}) H_{ij}^j \text{BR}(x_{-j})^{1-p} - \frac{1}{p+1} (1 - \|x_j\|_{p+1}^{p+1}) H_{ij}^j x_j^{1-p} \\
&= H_{ij}^j \left[ (\text{BR}(x_{-j}) - x_j) + \frac{1}{p+1} ((1 - \|\text{BR}(x_{-j})\|_{p+1}^{p+1}) \text{BR}(x_{-j})^{1-p} - (1 - \|x_j\|_{p+1}^{p+1}) x_j^{1-p}) \right].
\end{aligned}$$

*Deviation Incentive Gradient (Tsallis Entropy).* Finally, combining the derived terms gives:

$$\begin{aligned}
\nabla_{x_i} \mathcal{L}_{adi}(\mathbf{x}) &= -(\nabla_{x_i}^i - x_i^p \|\nabla_{x_i}^i\|_{1/p}) \\
&+ \sum_{j \neq i} H_{ij}^j \left[ (\text{BR}(x_{-j}) - x_j) + \frac{1}{p+1} ((1 - \|\text{BR}(x_{-j})\|_{p+1}^{p+1}) \text{BR}(x_{-j})^{1-p} - (1 - \|x_j\|_{p+1}^{p+1}) x_j^{1-p}) \right].
\end{aligned} \tag{35}$$

Note that in the limit of zero temperature, the gradient approaches

$$\nabla_{x_i} \mathcal{L}_{adi}(\mathbf{x}) \stackrel{p \rightarrow 0^+}{=} - \overbrace{(\nabla_{x_i}^i - 1 \|\nabla_{x_i}^i\|_\infty)}^{\text{policy gradient}} + \sum_{j \neq i} H_{ij}^j (\text{BR}_j - x_j). \tag{36}$$

The second component of the policy gradient term is orthogonal to the tangent space of the simplex, i.e., it does not contribute to movement along the simplex so it can be ignored in the limit of  $p \rightarrow 0^+$ .

Also, a Taylor series expansion of the adaptive Tsallis entropy around  $p = 0$  shows  $S_k^{\tau=p} = p s_k \mathcal{H}(x_k) + \mathcal{O}(p^2)$ , so the Tsallis entropy converges to a multiplicative constant of the Shannon entropy in the limit of zero entropy. If a similar homotopy exists for Tsallis entropy, maybe its limit point is the same limiting logit equilibrium as with Shannon entropy. We leave this to future research.

*Aside:* If you want to increase the entropy, just add a large constant to all payoffs which makes  $\text{BR} = \frac{1}{d}$  in the limit; it can be shown that  $\frac{1}{d}$  then becomes an equilibrium. Notice  $\text{BR}$  is invariant to multiplicative scaling of the payoffs. Therefore, deviation incentive is linear with respect to multiplicative scaling. One idea to decrease entropy is to subtract a constant from the payoffs such that they are still positive

but smaller. This can accomplish the desired effect, but will require more samples to estimate random variables with tiny values in their denominator. It seems like it won't be any more efficient than decreasing  $p$ .

### E.3 Shannon Entropy

The Nash equilibrium of utilities regularized with Shannon entropy is well known as the Quantal Response Equilibrium or Logit Equilibrium. The best response is a scaled softmax over the payoffs. We present the relevant intermediate gradients below.

$$S_k^\tau(x_k, x_{-k}) = -\tau \sum_i x_i \log(x_i) \quad (37)$$

$$\text{BR}(x_{-k}) = \text{softmax}\left(\frac{\nabla_{x_k}}{\tau}\right) \quad (38)$$

$$\nabla_{x_i} S_i^\tau(x_i, x_{-i}) = -\tau(\log(x_i) + 1) \quad (39)$$

$$\nabla_{x_i} S_j^\tau(x_j, x_{-j}) = 0 \quad (40)$$

$$J_{x_i}(\text{BR}(x_{-j})) = \frac{1}{\tau}(\text{diag}(\text{BR}_j) - \text{BR}_j \text{BR}_j^\top) H_{ji}^j \quad (41)$$

$$\nabla_{z_j} u_j^\tau(z_j, x_{-j})|_{\text{BR}_j, x_{-j}} = \nabla_{x_j}^j - \tau(\log(\text{BR}_j) + 1) \quad (42)$$

$$\nabla_{x_i} [u_i^\tau(\text{BR}(x_{-i}), x_{-i}) - u_i^\tau(x_i, x_{-i})] = -(\nabla_{x_i}^i - \tau(\log(\text{BR}_i) + 1)) \quad (43)$$

$$\nabla_{z_i} u_j^\tau(\text{BR}(x_{-j}), z_{-j})|_{\text{BR}_j, x_{-j}} = H_{ij}^j \text{BR}(x_{-j}) \quad (44)$$

$$\nabla_{z_i} u_j^\tau(x_j, z_{-j})|_{x_j, x_{-j}} = H_{ij}^j x_j \quad (45)$$

$$\begin{aligned} & \nabla_{x_i} [u_j^\tau(\text{BR}(x_{-j}), x_{-j}) - u_j^\tau(x_j, x_{-j})] \\ &= \left[ \frac{1}{\tau}(\text{diag}(\text{BR}_j) - \text{BR}_j \text{BR}_j^\top) H_{ji}^j \right]^\top (\nabla_{x_j}^j - \tau(\log(\text{BR}_j) + 1)) + H_{ij}^j \text{BR}(x_{-j}) - H_{ij}^j x_j \end{aligned} \quad (46)$$

*Deviation Incentive Gradient (Shannon Entropy).* Combining the derived terms gives:

$$\begin{aligned} \nabla_{x_i} \mathcal{L}_{adi}(\mathbf{x}) &= -(\nabla_{x_i}^i - \tau(\log(x_i) + 1)) \\ &+ \sum_{j \neq i} \left[ \frac{1}{\tau}(\text{diag}(\text{BR}_j) - \text{BR}_j \text{BR}_j^\top) H_{ji}^j \right]^\top (\nabla_{x_j}^j - \tau(\log(\text{BR}_j) + 1)) + H_{ij}^j [\text{BR}(x_{-j}) - x_j]. \end{aligned} \quad (47)$$

## F ABLATIONS

We introduce some additional notation here. A superscript indicates the temperature of the entropy regularizer, e.g.,  $\text{QRE}^{0.1}$  uses  $\tau = 0.1$  and  $\text{QRE}^{\text{auto}}$  anneals  $\tau$  as before. PED minimizes  $\mathcal{L}_{\text{adi}}$  without any entropy regularization or amortized estimates of payoff gradients (i.e., without the auxiliary variable  $y$ ).

### F.1 Bias re. §3.2+§4.1

Figure 7 demonstrates there exists a sweet spot for the amount of entropy regularization—too little and gradients are biased, too much and we solve for the Nash of a game we are not interested in.

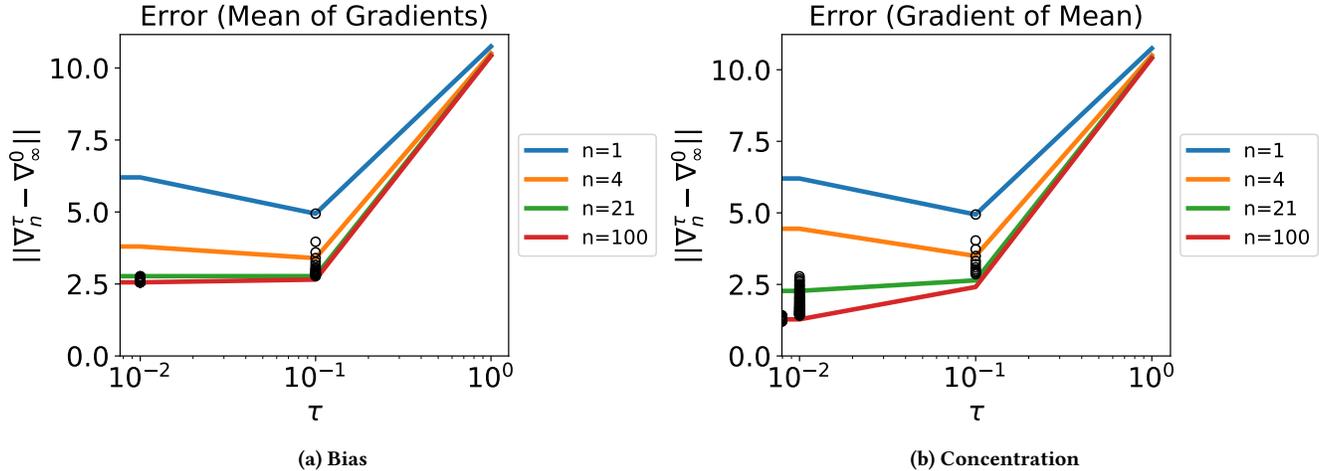


Figure 7: Bias-Bias Tradeoff on Blotto(10 coins, 3 fields, 4 players). Curves are drawn for samples sizes of  $n = \{1, 4, 21, 100\}$ . Circles denote the minimum of each curve for all  $n \in [1, 100]$ . Zero entropy regularization results in high gradient bias, i.e., stochastic gradients,  $\nabla_n^{\tau=0}$ , do not align well with the expected gradient,  $\nabla_\infty^{\tau=0}$ , where  $n$  is the number of samples. On the other hand, higher entropy regularization allows lower bias gradients but with respect to the entropy regularized utilities, not the unregularized utilities that we are interested in. The sweet spot lies somewhere in the middle. (a) SGD guarantees assume gradients are unbiased, i.e., the mean of sampled gradients is equal to the expected gradient in the limit of infinite samples  $n$ . Stochastic average deviation incentive gradients violate this assumption, the degree to which depends on the amount of entropy regularization  $\tau$  and number of samples  $n$ ;  $\tau = 10^{-2}$  appears to minimize the gradient bias for  $n = 100$  although with a nonzero asymptote around 2.5. (b) Computing a single stochastic gradient using more samples can reduce bias to zero in the limit. Note samples here refers to joint actions drawn from strategy profile  $x$ , not gradients as in (a). Additional samples makes gradient computation more expensive, but as we show later, these sample estimates can be amortized over iterations by reusing historical play. Both of the effects seen in (a) and (b) guide development of our proposed algorithm: (a) suggests using  $\tau > 0$  and (b) suggests reusing recent historical play to compute gradients (with  $\tau > 0$ ).

### F.2 Auxiliary $y$ re. §4.2

The introduction of auxiliary variables  $y_i$  are also supported by the results in Figure 8— $\text{QRE}^{0.0}$  is equivalent to PED and  ${}^y\text{QRE}^{0.0}$  is equivalent to PED augmented with  $y$ 's to estimate averages of payoff gradients.

### F.3 Annealing $\tau$ re. §3.1

ADIDAS includes temperature annealing, replacing the need to preset  $\tau$  with instead an ADI threshold  $\epsilon$ . Figure 9 compares this approach against other variants of the algorithm and shows this automated annealing mechanism reaches comparable final levels of ADI.

### F.4 Convergence re. §4.3

In Figure 9, FTRL and RM achieve low ADI quickly in some cases. FTRL has recently been proven not to converge to Nash, and this is suggested to be true of no-regret algorithms in general [30, 47]. Before proceeding, we demonstrate empirically in Figure 10 that FTRL and RM fail on games where minimizing  $\mathcal{L}_{\text{adi}}^\tau$  still makes progress, even without an annealing schedule.

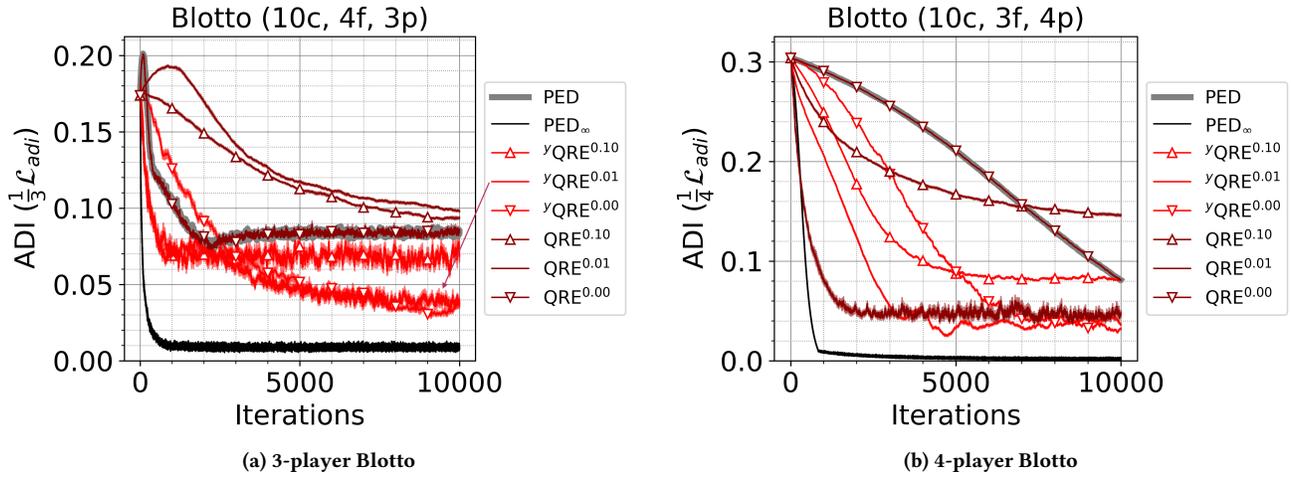


Figure 8: Adding an appropriate level of entropy can accelerate convergence (compare PED to QRE $^{0.01}$  in (b)). And amortizing estimates of joint play using  $y$  can reduce gradient bias, further improving performance (e.g., compare QRE $^{0.00}$  to  $y$ QRE $^{0.00}$  in (a) or (b)).

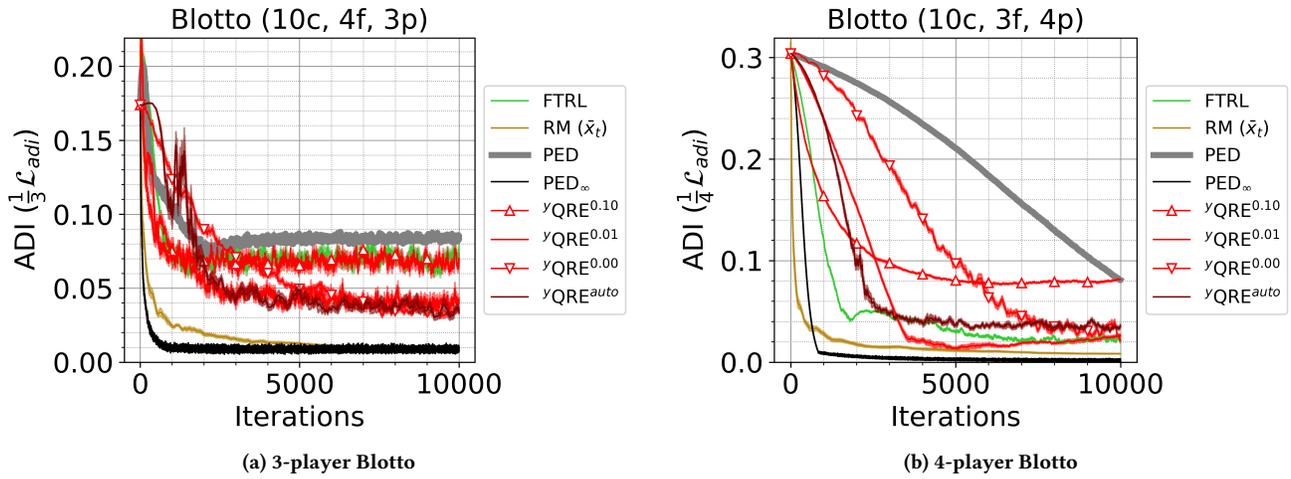


Figure 9: Average deviation incentive of the symmetric joint strategy  $x^{(t)}$  is plotted against algorithm iteration  $t$ . Despite FTRL’s lack of convergence guarantees, it converges quickly in these games.

## F.5 ADI stochastic estimate

Computing ADI exactly requires the full payoff tensor, so in very large games, we must estimate ADI. Figure 11 shows how estimates of  $\mathcal{L}_{adi}$  computed from historical play track their true expected value throughout training.

## G EXPERIMENTS REPEATED WITH ATE

*Bias re. §3.2+§4.1.* We first empirically verify that adding an entropy regularizer to the player utilities introduces a trade-off: set entropy regularization too low and the best-response operator will have high bias; set entropy regularization too high and risk solving for the Nash of a game we are not interested in. Figure 12 shows there exists a sweet spot in the middle for moderate amounts of regularization (temperatures).

*Auxiliary  $y$  re. §4.2.* The introduction of auxiliary variables  $y_i$  are supported by the results in Figure 13—ATE $^{0.0}$  is equivalent to PED and  $y$ ATE $^{0.0}$  is equivalent to PED augmented with  $y$ ’s to estimate averages of payoff gradients.

In Figure 13, we also see a more general relationship between temperature and convergence rate. Higher temperatures appear to result in faster initial convergence ( $\mathcal{L}_{adi}$  spikes initially in Figure 13a for  $\tau < 0.1$ ) and lower variance but higher asymptotes, while the opposite holds for lower temperatures. These results suggest annealing the temperature over time to achieve fast initial convergence and lower asymptotes.

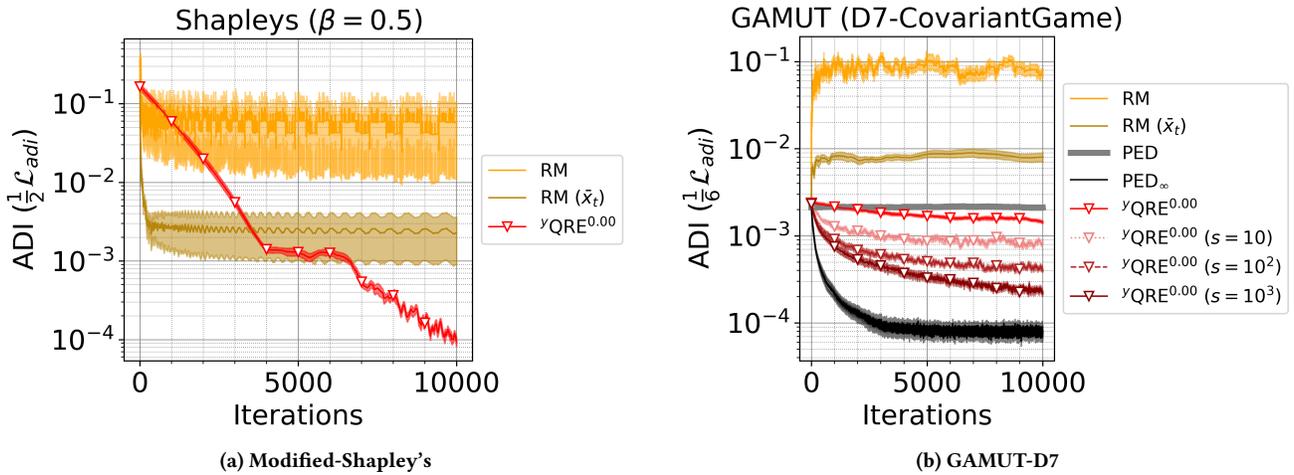


Figure 10: ADIDAS reduces  $\mathcal{L}_{adi}$  in both games. In game (a), created by Ostrovski and van Strien [51], to better test performance,  $x$  is initialized randomly rather than with the uniform distribution because the Nash is at uniform. In (b), computing gradients using full expectations (in black) results in very low ADI. Computing gradients using only single samples plus historical play allows a small reduction in ADI. More samples (e.g.,  $n = 10^3$ ) allows further reduction.

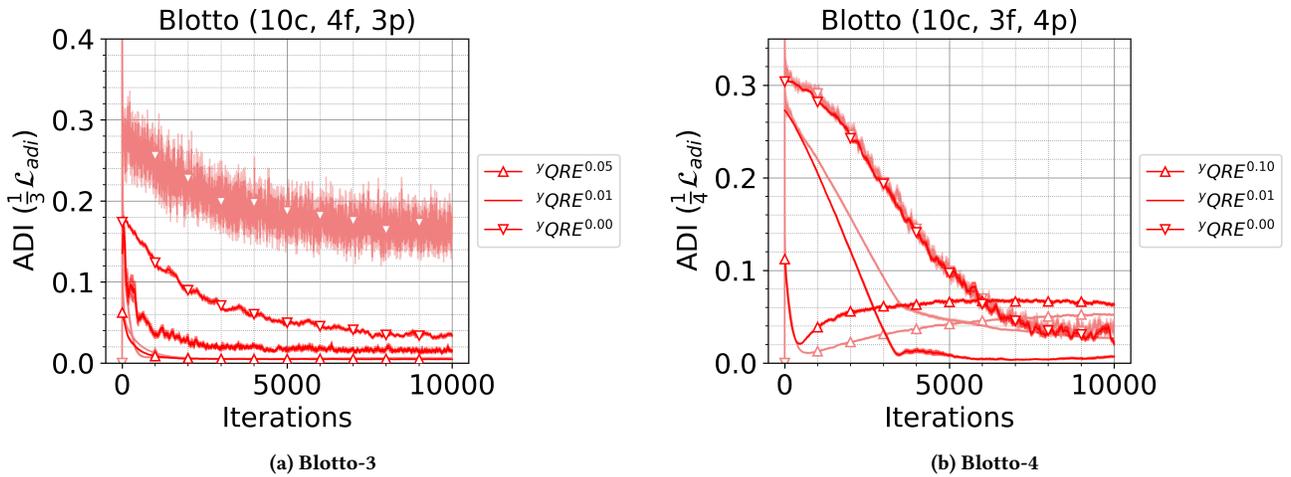


Figure 11: Accuracy of running estimate of  $\mathcal{L}_{adi}^\tau$  computed from  $y^{(t)}$  (in light coral) versus true value (in red).

Lower variance should also be possible by carefully annealing the learning rate to allow  $y$  to accurately perform tracking. Fixed learning rates were used here; we leave investigating learning rate schedules to future work.

Figure 13b shows how higher temperatures (through a reduction in gradient bias) can result in accelerated convergence.

*Annealing  $\tau$  re. §3.1.* ADIDAS includes temperature annealing replacing the need for setting the hyperparameter  $\tau$  with instead an ADI threshold  $\epsilon$ . Figure 14 compares this approach against several other variants of the algorithm and shows this automated annealing mechanism reaches comparable final levels of ADI.

*Convergence re. §4.3.* In Figure 14, FTRL and RM achieve low levels of ADI quickly in some cases. FTRL has recently been proven not to converge to Nash, and this is suggested to be true of no-regret algorithms such as RM in general [30, 47]. Before proceeding, we demonstrate empirically in Figure 15 that FTRL and RM fail on some games where ADIDAS still makes progress.

*Large-Scale re §4.4.* Computing ADI exactly requires the full payoff tensor, so in very large games, we must estimate the ADI. Figure 16 shows how estimates of  $\mathcal{L}_{adi}$  computed from historical play track their true expected value throughout training.

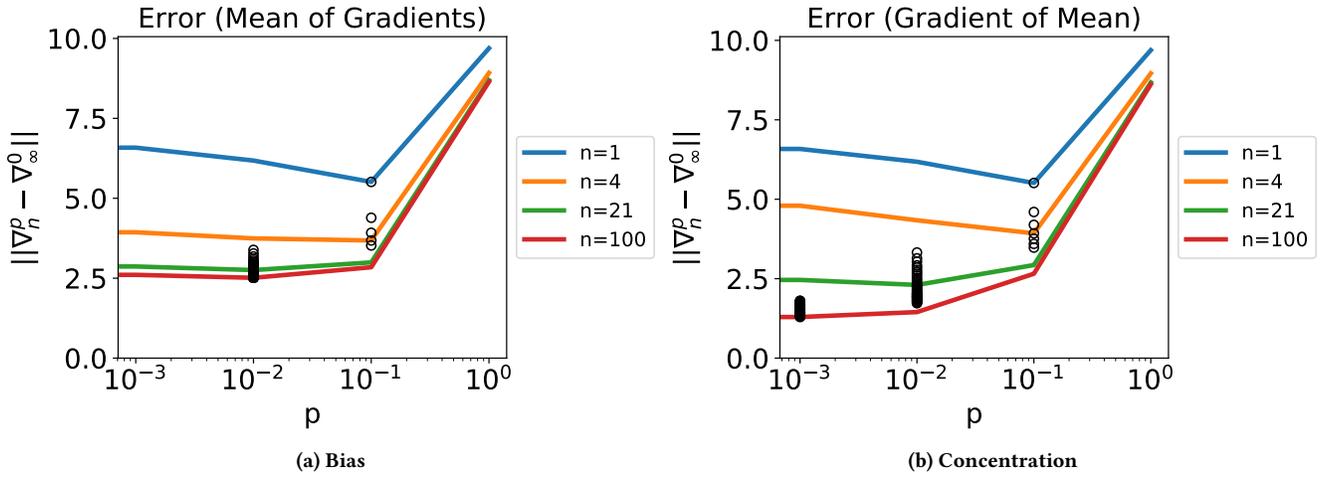


Figure 12: Bias-Bias Tradeoff on Blotto(10 coins, 3 fields, 4 players). Curves are drawn for samples sizes of  $n = \{1, 4, 21, 100\}$ . Circles denote the minimum of each curve for all  $n \in [1, 100]$ . Zero entropy regularization results in high gradient bias, i.e., stochastic gradients,  $\nabla_n^{\tau=0}$ , do not align well with the expected gradient,  $\nabla_\infty^{\tau=0}$ , where  $n$  is the number of samples. On the other hand, higher entropy regularization allows lower bias gradients but with respect to the entropy regularized utilities, not the unregularized utilities that we are interested in. The sweet spot lies somewhere in the middle. (a) SGD guarantees assume gradients are unbiased, i.e., the mean of sampled gradients is equal to the expected gradient in the limit of infinite samples  $n$ . Stochastic average deviation incentive gradients violate this assumption, the degree to which depends on the amount of entropy regularization  $\tau$  and number of samples  $n$ ;  $p = 10^{-2}$  appears to minimize the gradient bias for  $n = 100$  although with a nonzero asymptote around 2.5. (b) Computing a single stochastic gradient using more samples can reduce bias to zero in the limit. Note samples here refers to joint actions from strategy profile  $x$ , not gradients as in (a). Additional samples makes gradient computation more expensive, but as we show later, these sample estimates can be amortized over iterations by reusing historical play. Both the effects seen in (a) and (b) guide development of our proposed algorithm: (a) suggests using  $\tau > 0$  and (b) suggests reusing recent historical play to compute gradients (with  $\tau > 0$ ).

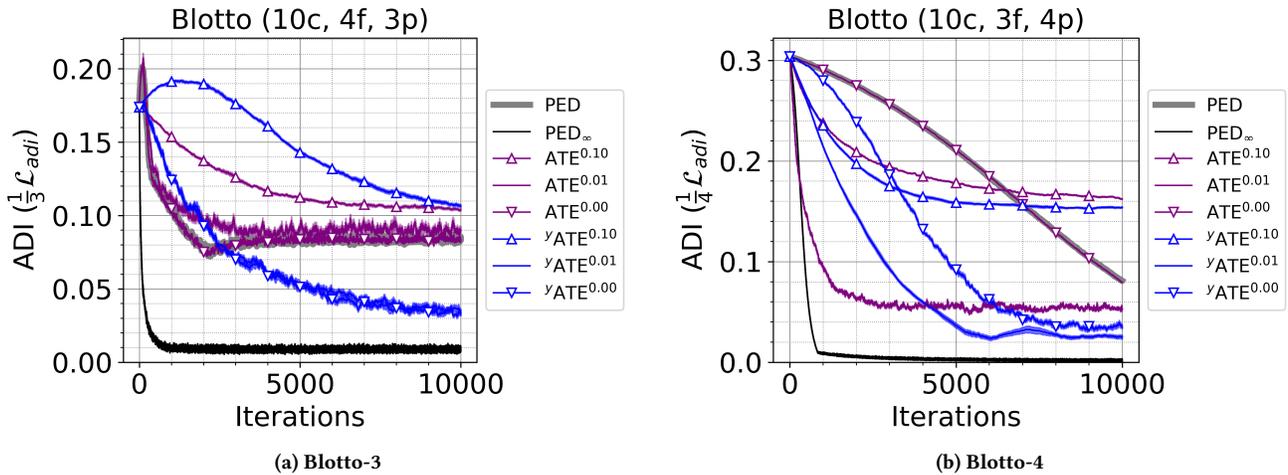


Figure 13: (a) 3-player Blotto game; (b) 4-player Blotto game. Adding an appropriate level of entropy (e.g.,  $\tau = 0.01$ ) can accelerate convergence (compare PED to  $\text{ATE}^{0.01}$  in (b)). And amortizing estimates of joint play can reduce gradient bias, further improving performance (e.g., compare  $\text{ATE}^{0.01}$  to  $\text{yATE}^{0.01}$  in (a) or (b)).

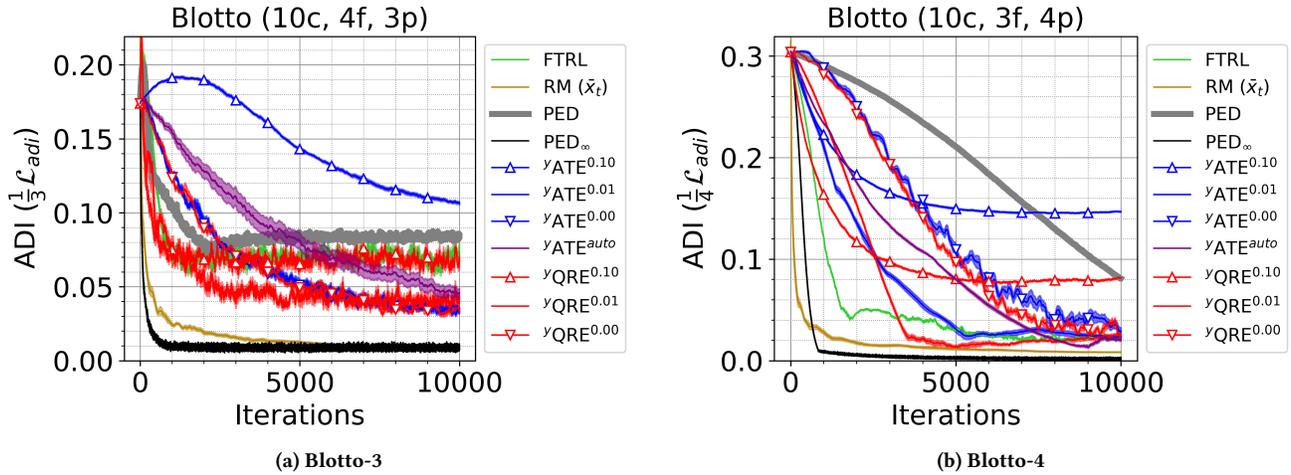


Figure 14: (a) 3-player Blotto game; (b) 4-player Blotto game. The maximum a single agent can exploit the symmetric joint strategy  $x^{(t)}$  is plotted against algorithm iteration  $t$ . Despite FTRL’s lack of convergence guarantees, it converges quickly in these Blotto games.

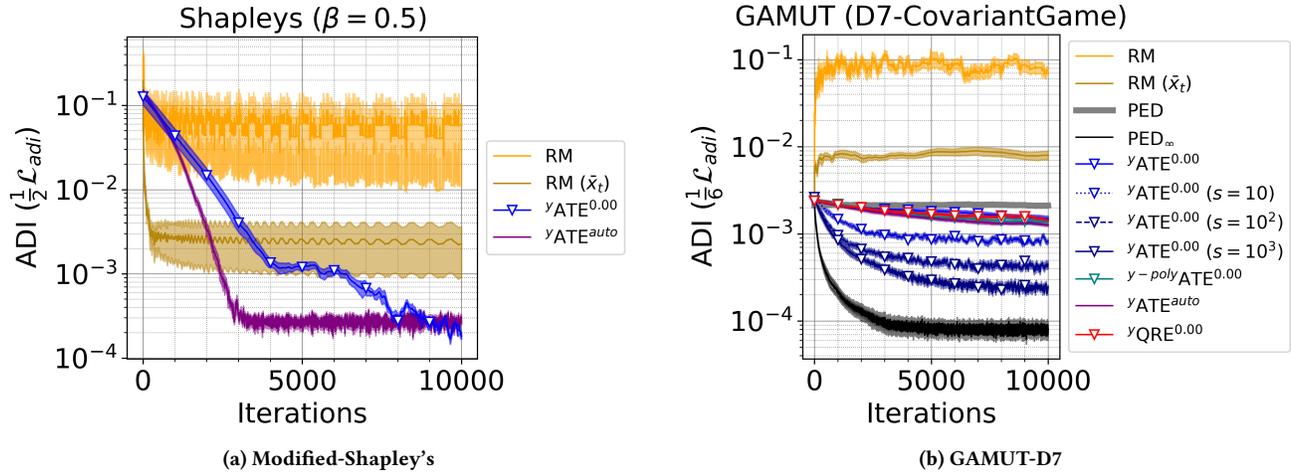


Figure 15: (a) Modified-Shapley’s; (b) GAMUT-D7. Deviation incentive descent reduces  $\mathcal{L}_{adi}$  in both games. In (a), to better test the performance of the algorithms,  $x$  is initialized randomly rather than with the uniform distribution because the Nash is at uniform. In (b), computing ADI gradients using full expectations (in black) results in very low levels of ADI. Computing estimates using only single samples plus historical play allows a small reduction in ADI. More samples (e.g.,  $n = 10^3$ ) allows further reduction.

## H COMPARISON AGAINST ADDITIONAL ALGORITHMS

### H.1 ED and FP Fail

We chose not to include Exploitability Descent (ED) or Fictitious Play (FP) in the main body as we considered them to be “straw men”. ED is only expected to converge in 2-player, zero-sum games. FP is non-convergent in some 2-player games as well [32]. We run ED and FP with true expected gradients & best responses ( $s=\infty$ ) on the 3 player game in Figure 17 to convince the reader that failure to converge is not due to stochasticity.

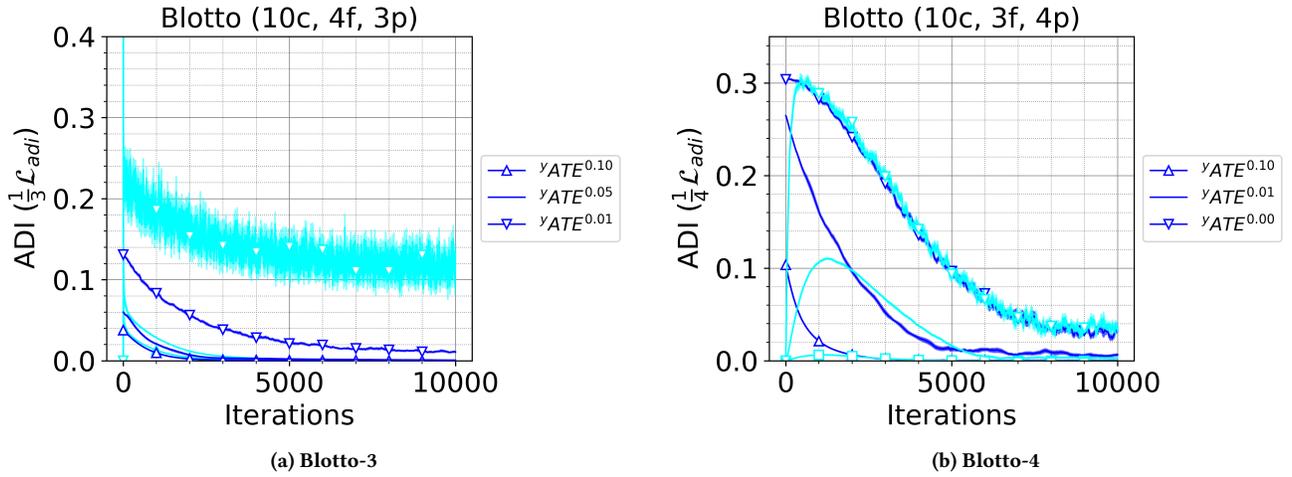


Figure 16: Accuracy of running estimate of  $\mathcal{L}_{adi}^\tau$  computed from  $y^{(t)}$  (in light coral) versus true value (in blue).

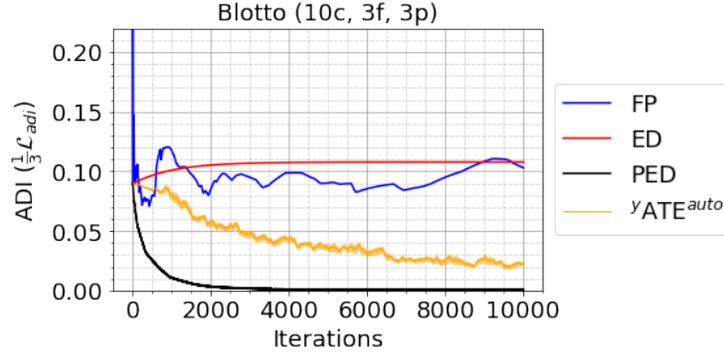


Figure 17: FP, ED, PED access the full tensor.  $y^{ATE^{auto}}$  samples.

## H.2 Gambit Solvers

We ran all applicable gambit solvers on the 4-player, 10-coin, 3-field Blotto game (command listed below). All solvers fail to return a Nash equilibrium except gambit-enumpoly which returns all 36 permutations of the following pure, non-symmetric Nash equilibrium:

$$x^* = [(10, 0, 0), (10, 0, 0), (0, 10, 0), (0, 0, 10)] \quad (48)$$

where each of the four players places 10 coins on one of the three fields.

- gambit-enumpoly
- gambit-gnm
- gambit-ipa
- gambit-liap
- gambit-simpdiv
- gambit-logit

Command:

```
timeout 3600s gambit-enumpoly -H < blotto_10_3_4.nfg >> enumpoly.txt; timeout 3600s gambit-gnm < blotto_10_3_4.nfg >>
gnm.txt; timeout 3600s gambit-ipa < blotto_10_3_4.nfg >> ipa.txt; timeout 3600s gambit-liap < blotto_10_3_4.nfg >>
liap.txt; timeout 3600s gambit-simpdiv < blotto_10_3_4.nfg >> simpdiv.txt; timeout 3600s gambit-logit -m 1.0 -e <
blotto_10_3_4.nfg >> logit.txt
```

# I ADDITIONAL GAME DOMAINS

## I.1 Diplomacy Experiments - Subsampled Games

Figure 18 runs a comparison on 40 subsampled tensors (7-players, 4-actions each) taken from the 40 turns of a single Diplomacy match. The four actions selected for each player are sampled from the corresponding player’s trained policy.

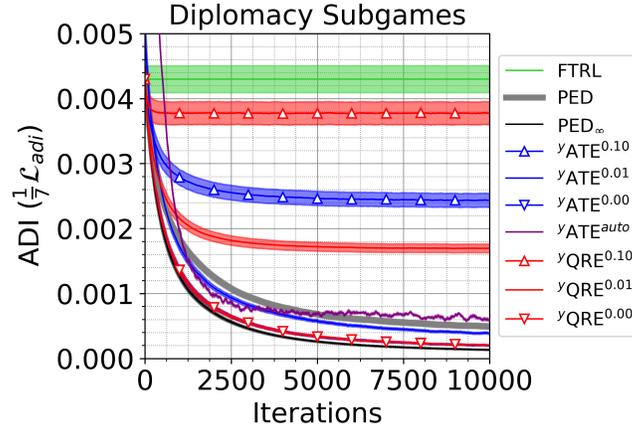


Figure 18: Subsampled games.

Figure 19 runs a comparison on two Diplomacy meta-games, one with 5 bots trained using Fictitious Play and the other with bots trained using Iterated Best Response (IBR) –these are the same meta-games analyzed in Figure 3 of [2].

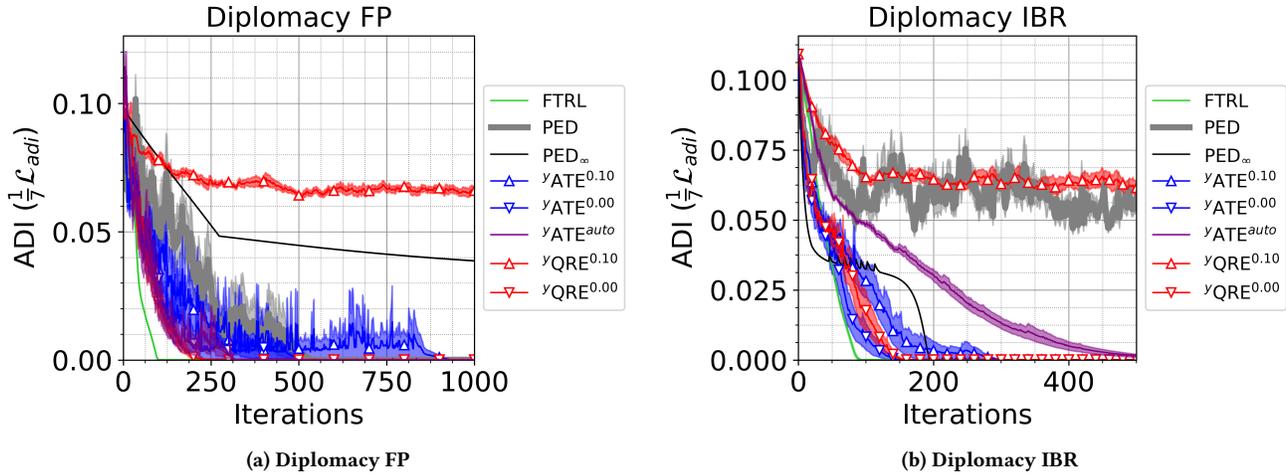


Figure 19: (a) FP; (b) IBR. The maximum a single agent can exploit the symmetric joint strategy  $x^{(t)}$  is plotted against algorithm iteration  $t$ . Many of the algorithms quickly achieve near zero  $\mathcal{L}_{adi}$ , so unlike in the other experiments, hyperparameters are selected according according to the earliest point at which exploitability falls below 0.01 with ties split according to the final value.

Figure 5 demonstrates an empirical game theoretic analysis [36, 61, 62] of a large symmetric 7-player Diplomacy meta-game where each player elects 1 of 5 trained bots to play on their behalf. In this case, the expected value of each entry in the payoff tensor represents a winrate. Each entry can only be estimated by simulating game play, and the result of each game is a Bernoulli random variable. To obtain a winrate estimate within 0.01 of the true estimate with probability 95%, a Chebyshev bound implies more than 223 samples are needed. The symmetric payoff tensor contains 330 unique entries, requiring over 74 thousand games in total. In the experiment below, ADIDAS achieves negligible ADI in less than 7 thousand iterations with 50 samples of joint play per iteration ( $\approx 5\times$  the size of the tensor).

## I.2 Diplomacy Experiments - Empirical Game Theoretic Analysis

Figure 20 repeats the computation of Figure 6 with a smaller auxiliary learning rate  $\eta_y$  and achieves better results.

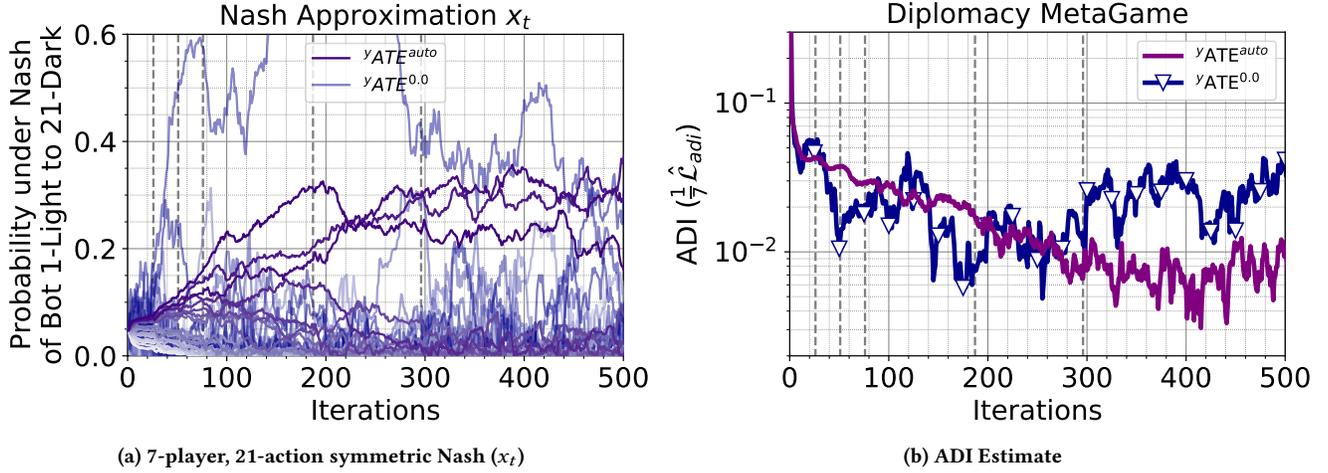


Figure 20: (a) Evolution of the symmetric Nash approximation; (b) ADI estimated from auxiliary variable  $y_t$ . Black vertical lines indicate the temperature  $\tau$  was annealed. Auxiliary learning rate  $\eta_y = 1/25$ . In addition to the change in  $\eta_y$  from  $1/10$ , also note the change in axes limits versus Figure 6.

## I.3 El Farol Bar Stage Game

We compare ADIDAS variants and regret matching in Figure 21 on the 10-player symmetric El Farol Bar stage game with hyperparameters  $n = 10$ ,  $c = 0.7$ ,  $C = nc$ ,  $B = 0$ ,  $S = 1$ ,  $G = 2$  (see Section 3.1, The El Farol stage game in [63]). Recall that the homotopy that ADIDAS attempts to trace is displayed in Figure 2b of the main body.

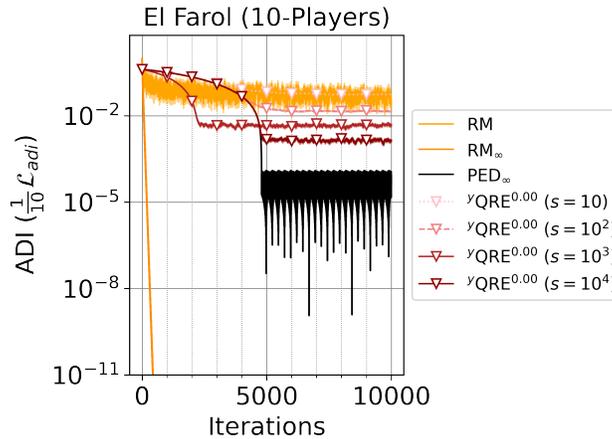


Figure 21: (10-player, 2-action El Farol stage game) ADIDAS and regret matching both benefit from additional samples. Both find the same unique mixed Nash equilibrium. In this case, regret matching finds it much more quickly.

## J DESCRIPTION OF DOMAINS

### J.1 Modified Shapley's

The modified Shapley's game mentioned in the main body (Figure 4a) is defined in Table 5 [51].

### J.2 Colonel Blotto

Despite its apparently simplicity, Colonel Blotto is a complex challenge domain and one under intense research [1, 7–9, 13].

1	0	$\beta$
$\beta$	1	0
0	$\beta$	1

(a) Player A's Payoff Matrix

$-\beta$	1	0
0	$-\beta$	1
1	0	$-\beta$

(b) Player B's Payoff Matrix

Table 5: (a) Player A; (b) Player B. We set  $\beta = 0.5$  in experiments and subtract  $-\beta$  from each payoff matrix to ensure payoffs are non-negative; ATE requires non-negative payoffs.

## K CONNECTIONS TO OTHER ALGORITHMS

### K.1 Consensus Algorithm

ADIDAS with Tsallis entropy and temperature fixed to  $\tau = p = 1$  recovers the regularizer proposed for the Consensus algorithm [48] plus entropy regularization. To see this, recall from Appx. E.2, the Tsallis entropy regularizer:

$$S_k^{\tau=p=1}(x_k, x_{-k}) = \frac{s_k}{2} \left(1 - \sum_m x_{km}^2\right) \quad (49)$$

where  $s_k = \left(\sum_m (\nabla_{x_{km}}^k)^{1/1}\right)^1 = \|\nabla_{x_k}^k\|_{1/1}$  is treated as a constant w.r.t.  $\mathbf{x}$ .

In the case with  $\tau = p = 1$ ,  $\text{BR}_k = \text{BR}(x_{-k}) = \frac{1}{s_k} \nabla_{x_k}^k$  where we have assumed the game has been offset by a constant so that it contains only positive payoffs. Plugging these into the definition of  $\mathcal{L}_{adi}^\tau$ , we find

$$\mathcal{L}_{adi}^\tau(\mathbf{x}) = \sum_k u_k^\tau(\text{BR}_k, x_{-k}) - u_k^\tau(x_k, x_{-k}) \quad (50)$$

$$\approx \sum_k u_k \left( \frac{1}{s_k} \nabla_{x_k}^k, x_{-k} \right) - u_k(x_k, x_{-k}) \quad (51)$$

$$= \sum_k \frac{1}{s_k} \underbrace{\|\nabla_{x_k}^k\|^2}_{\text{consensus regularizer}} - x_k^\top \nabla_{x_k}^k. \quad (52)$$

Note the Consensus regularizer can also be arrived at by replacing the best response with a 1-step gradient ascent lookahead, i.e.,  $\text{BR}_k = x_k + \eta \nabla_{x_k}^k$ :

$$\mathcal{L}_{adi}^\tau(\mathbf{x}) = \sum_k u_k^\tau(\text{BR}_k, x_{-k}) - u_k^\tau(x_k, x_{-k}) \quad (53)$$

$$\approx \sum_k u_k(x_k + \eta \nabla_{x_k}^k, x_{-k}) - u_k(x_k, x_{-k}) \quad (54)$$

$$= \sum_k (x_k^\top \nabla_{x_k}^k) + \eta \|\nabla_{x_k}^k\|^2 - (x_k^\top \nabla_{x_k}^k) \quad (55)$$

$$= \sum_k \eta \|\nabla_{x_k}^k\|^2. \quad (56)$$

### K.2 Exploitability Descent as Extragradient

In normal-form games, Exploitability Descent (ED) [42] is equivalent to Extragradient [38] (or Mirror Prox [37]) with an infinite intermediate step size. Recall  $\text{BR}_k = \arg \max_{x \in \Delta^{m_k-1}} u_k(x_k, x_{-k}) = \arg \max_{x \in \Delta^{m_k-1}} x^\top \nabla_{x_k}^k$ . Using the convention that ties between actions result in vectors that distribute a 1 uniformly over the maximizers, the best response can be rewritten as  $\text{BR}_k = \lim_{\hat{\eta} \rightarrow \infty} \Pi[x_k + \hat{\eta} \nabla_{x_k}^k]$  where  $\Pi$  is the Euclidean projection onto the simplex. Define  $F(\mathbf{x})$  such that its  $k$ th component  $F(\mathbf{x})_k = -\nabla_{x_k}^k = -\nabla_{x_k}^k(x_{-k})$  where we have simply introduced  $x_{-k}$  in parentheses to emphasize that player  $k$ 's gradient is a function of  $x_{-k}$  only, and not  $x_k$ . Equations without subscripts imply they are applied in parallel over the players.

ED executes the following update in parallel for all players  $k$ :

$$x_{k+1} \leftarrow \Pi[x_k + \eta \nabla_{x_k} \{u_k(x_k, x_{-k})\}|_{x_{-k}=\text{BR}_{-k}}]. \quad (57)$$

Define  $\hat{x}_k = x_k - \hat{\eta}F(x)_k$ . And as an abuse of notation, let  $\hat{x}_{-k} = x_{-k} - \hat{\eta}F(x)_{-k}$ . Extragradient executes the same update in parallel for all players  $k$ :

$$x_{k+1} \leftarrow \Pi[x_k - \eta F(\Pi[x - \hat{\eta}F(x)])_k] \quad (58)$$

$$= \Pi[x_k - \eta F(\Pi[x + \hat{\eta}\nabla_x])_k] \quad (59)$$

$$= \Pi[x_k - \eta F(\text{BR})_k] \quad (60)$$

$$= \Pi[x_k + \eta \nabla_{x_k} \{u_k(x_k, x_{-k})\}|_{x_{-k}=\text{BR}_{-k}}]. \quad (61)$$

Extragradient is known to converge in two-player zero-sum normal form games given an appropriate step size scheme. The main property that Extragradient relies on is monotonicity of the vector function  $F$ . All two-player zero-sum games induce a monotone  $F$ , however, this is not true in general of two-player general-sum or games with more players. ED is only proven to converge for two-player zero-sum, but this additional connection provides an additional reason why we do not expect ED to solve many-player general-sum normal-form games, which are the focus of this work. Please see Appx. H.1 for an experimental demonstration.

## L PYTHON CODE

For the sake of reproducibility we have included code in python+numpy.

```

1 """
2 Copyright 2020 ADIDAS Authors.
3
4
5 Licensed under the Apache License, Version 2.0 (the "License");
6 you may not use this file except in compliance with the License.
7 You may obtain a copy of the License at
8
9 https://www.apache.org/licenses/LICENSE-2.0
10
11 Unless required by applicable law or agreed to in writing, software
12 distributed under the License is distributed on an "AS IS" BASIS,
13 WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
14 See the License for the specific language governing permissions and
15 limitations under the License.
16 """
17 import numpy as np
18 from scipy import special
19
20 def simplex_project_grad(g):
21     """Project gradient onto tangent space of simplex."""
22     return g - g.sum() / g.size

```

Listing 1: Header.

```

1 def gradients_qre_nonsym(dist, y, anneal_steps, payoff_matrices,
2                          num_players, temp=0., proj_grad=True,
3                          exp_thresh=1e-3, lrs=(1e-2, 1e-2),
4                          logit_clip=-1e5):
5     """Computes exploitability gradient and aux variable gradients.
6
7     Args:
8         dist: list of 1-d np.arrays, current estimate of nash
9         y: list of 1-d np.arrays, current est. of payoff gradient
10        anneal_steps: int, elapsed num steps since last anneal
11        payoff_matrices: dict with keys as tuples of agents (i, j) and
12        values of (2 x A x A) arrays, payoffs for each joint action.
13        keys are sorted and arrays are indexed in the same order.
14        num_players: int, number of players
15        temp: non-negative float, default 0.
16        proj_grad: bool, if True, projects dist gradient onto simplex
17        exp_thresh: ADI threshold at which temp is annealed
18        lrs: tuple of learning rates (lr_x, lr_y)
19        logit_clip: float, minimum allowable logit
20
21     Returns:
22        gradient of ADI w.r.t. (dist, y, anneal_steps)
23        temperature (possibly annealed, i.e., reduced)
24        unregularized ADI (stochastic estimate)
25        shannon regularized ADI (stochastic estimate)
26     """

```

```

26 # first compute policy gradients and player effects (fx)
27 policy_gradient = []
28 other_player_fx = []
29 grad_y = []
30 unreg_exp = []
31 reg_exp = []
32 for i in range(num_players):
33
34     nabla_i = np.zeros_like(dist[i])
35     for j in range(num_players):
36         if j == i:
37             continue
38         if i < j:
39             hess_i_ij = payoff_matrices[(i, j)][0]
40         else:
41             hess_i_ij = payoff_matrices[(j, i)][1].T
42
43     nabla_ij = hess_i_ij.dot(dist[j])
44     nabla_i += nabla_ij / float(num_players - 1)
45
46     grad_y.append(y[i] - nabla_i)
47
48     if temp >= 1e-3: # numerical under/overflow for temp < 1e-3
49         br_i = special.softmax(y[i] / temp)
50         br_i_mat = (np.diag(br_i) - np.outer(br_i, br_i)) / temp
51         log_br_i_safe = np.clip(np.log(br_i), logit_clip, 0)
52         br_i_policy_gradient = nabla_i - temp * (log_br_i_safe + 1)
53     else:
54         power = np.inf
55         s_i = np.linalg.norm(y[i], ord=power)
56         br_i = np.zeros_like(dist[i])
57         maxima_i = (y[i] == s_i)
58         br_i[maxima_i] = 1. / maxima_i.sum()
59         br_i_mat = np.zeros((br_i.size, br_i.size))
60         br_i_policy_gradient = np.zeros_like(br_i)
61
62     policy_gradient_i = np.array(nabla_i)
63     if temp > 0:
64         log_dist_i_safe = np.clip(np.log(dist[i]), logit_clip, 0)
65         policy_gradient_i -= temp * (log_dist_i_safe + 1)
66     policy_gradient.append(policy_gradient_i)
67
68     unreg_exp_i = np.max(y[i]) - y[i].dot(dist[i])
69     unreg_exp.append(unreg_exp_i)
70
71     entr_br_i = temp * special.entr(br_i).sum()
72     entr_dist_i = temp * special.entr(dist[i]).sum()
73
74     reg_exp_i = y[i].dot(br_i - dist[i]) + entr_br_i - entr_dist_i
75     reg_exp.append(reg_exp_i)
76
77     other_player_fx_i = (br_i - dist[i])
78     other_player_fx_i += br_i_mat.dot(br_i_policy_gradient)
79     other_player_fx.append(other_player_fx_i)
80
81 # then construct ADI gradient
82 grad_dist = []
83 for i in range(num_players):
84
85     grad_dist_i = -policy_gradient[i]
86     for j in range(num_players):
87         if j == i:
88             continue
89         if i < j:
90             hess_j_ij = payoff_matrices[(i, j)][1]
91         else:
92             hess_j_ij = payoff_matrices[(j, i)][0].T
93
94     grad_dist_i += hess_j_ij.dot(other_player_fx[j])

```

```

95
96     if proj_grad:
97         grad_dist_i = simplex_project_grad(grad_dist_i)
98
99         grad_dist.append(grad_dist_i)
100
101     unreg_exp_mean = np.mean(unreg_exp)
102     reg_exp_mean = np.mean(reg_exp)
103
104     _, lr_y = lrs
105     if (reg_exp_mean < exp_thresh) and (anneal_steps >= 1 / lr_y):
106         temp = np.clip(temp / 2., 0., 1.)
107         if temp < 1e-3: # consistent with numerical issue above
108             temp = 0.
109         grad_anneal_steps = -anneal_steps
110     else:
111         grad_anneal_steps = 1
112
113     return ((grad_dist, grad_y, grad_anneal_steps), temp,
114            unreg_exp_mean, reg_exp_mean)

```

**Listing 2: ADIDAS Gradient.**

```

1 def gradients_ate_sym(dist, y, anneal_steps, payoff_matrices,
2                       num_players, p=1, proj_grad=True,
3                       exp_thresh=1e-3, lrs=(1e-2, 1e-2)):
4     """Computes ADI gradient and aux variable gradients.
5
6     Args:
7         dist: list of 1-d np.arrays, current estimate of nash
8         y: list of 1-d np.arrays, current est. of payoff gradient
9         anneal_steps: int, elapsed num steps since last anneal
10        payoff_matrices: dict with keys as tuples of agents (i, j) and
11            values of (2 x A x A) arrays, payoffs for each joint action.
12            keys are sorted and arrays are indexed in the same order.
13        num_players: int, number of players
14        p: float in [0, 1], Tsallis entropy-regularization
15        proj_grad: bool, if True, projects dist gradient onto simplex
16        exp_thresh: ADI threshold at which p is annealed
17        lrs: tuple of learning rates (lr_x, lr_y)
18
19    Returns:
20        gradient of ADI w.r.t. (dist, y, anneal_steps)
21        temperature, p (possibly annealed, i.e., reduced)
22        unregularized ADI (stochastic estimate)
23        tsallis regularized ADI (stochastic estimate)
24    """
25    nabla = payoff_matrices[0].dot(dist)
26    if p >= 1e-2: # numerical under/overflow when power > 100.
27        power = 1. / float(p)
28        s = np.linalg.norm(y, ord=power)
29        if s == 0:
30            br = np.ones_like(y) / float(y.size) # uniform dist
31        else:
32            br = (y / s)**power
33    else:
34        power = np.inf
35        s = np.linalg.norm(y, ord=power)
36        br = np.zeros_like(dist)
37        maxima = (y == s)
38        br[maxima] = 1. / maxima.sum()
39
40    unreg_exp = np.max(y) - y.dot(dist)
41    br_inv_sparse = 1 - np.sum(br**(p + 1))
42    dist_inv_sparse = 1 - np.sum(dist**(p + 1))
43    entr_br = s / (p + 1) * br_inv_sparse
44    entr_dist = s / (p + 1) * dist_inv_sparse
45    reg_exp = y.dot(br - dist) + entr_br - entr_dist
46
47    entr_br_vec = br_inv_sparse * br**(1 - p)

```

```

47  entr_dist_vec = dist_inv_sparse * dist**(1 - p)
48
49  policy_gradient = nabla - s * dist**p
50  other_player_fx = (br - dist)
51  other_player_fx += 1 / (p + 1) * (entr_br_vec - entr_dist_vec)
52
53  other_player_fx_translated = payoff_matrices[1].dot(
54      other_player_fx)
55  grad_dist = -policy_gradient
56  grad_dist += (num_players - 1) * other_player_fx_translated
57  if proj_grad:
58      grad_dist = simplex_project_grad(grad_dist)
59  grad_y = y - nabla
60
61  _, lr_y = lrs
62  if (reg_exp < exp_thresh) and (anneal_steps >= 1 / lr_y):
63      p = np.clip(p / 2., 0., 1.)
64      if p < 1e-2: # consistent with numerical issue above
65          p = 0.
66      grad_anneal_steps = -anneal_steps
67  else:
68      grad_anneal_steps = 1
69
70  return ((grad_dist, grad_y, grad_anneal_steps), p, unreg_exp,
71      reg_exp)

```

**Listing 3: ADIDAS Gradient (assuming symmetric Nash and with Tsallis entropy).**