

# Small Nash Equilibrium Certificates in Very Large Games

Brian Hu Zhang,<sup>1</sup> Tuomas Sandholm<sup>1,2,3,4</sup>

<sup>1</sup> Computer Science Department, Carnegie Mellon University

<sup>2</sup> Strategic Machine, Inc.

<sup>3</sup> Strategy Robot, Inc.

<sup>4</sup> Optimized Markets, Inc.

## Abstract

In many game settings, the game is not explicitly given but is only accessible by playing it. While there have been impressive demonstrations in such settings, prior techniques have not offered safety guarantees, that is, guarantees on the game-theoretic exploitability of the computed strategies. In this paper we introduce an approach that shows that it is possible to provide exploitability guarantees in such settings without ever exploring the entire game. We introduce a notion of a *certificate* of an extensive-form approximate Nash equilibrium. For verifying a certificate, we give an algorithm that runs in time linear in the size of the certificate rather than the size of the whole game. In zero-sum games, we further show that an optimal certificate—given the exploration so far—can be computed with any standard game-solving algorithm (e.g., using a linear program or counterfactual regret minimization). However, unlike in the cases of normal form or perfect information, we show that certain families of extensive-form games do not have small approximate certificates, even after making extremely nice assumptions on the structure of the game. Despite this difficulty, we find experimentally that very small certificates, even exact ones, often exist in large and even in infinite games. Overall, our approach enables one to try one’s favorite exploration strategies while offering exploitability guarantees, thereby decoupling the exploration strategy from the equilibrium-finding process.

## 1 Introduction

Recent years have witnessed AI breakthroughs in games such as poker (Bowling et al. 2015; Moravčík et al. 2017; Brown and Sandholm 2017b, 2019b) where the rules are given. In many important applications—such as many war games and finance simulations—the rules are only given via *black-box* access, that is, via playing the game (Wellman 2006; Lanctot et al. 2017), and one can try to construct good strategies by self play. In such settings, deep reinforcement learning techniques are typically used today (Heinrich and Silver 2016; Silver et al. 2016; Lanctot et al. 2017; Srinivasan et al. 2018; Vinyals et al. 2019; Berner et al. 2019). However, such methods lack the guarantee of low (or zero) *exploitability* that game-theoretic solving techniques offer.

Prior to our paper, to compute exploitability of a strategy, one needed to compute the other player’s best response

to it, which relies on the game being known. Sampling approaches to equilibrium finding have been suggested, but their regret guarantees are vacuous unless the algorithms touch at least as many information sets as there are in the game (Lanctot et al. 2009; Srinivasan et al. 2018; Zhou, Li, and Zhu 2020). A recent PAC-learning algorithm has logarithmic sample complexity for *pure* maxmin strategies in *normal-form* games; it extends to some infinite games, but not effectively to mixed strategies in extensive-form games Marchesi, Trovò, and Gatti (2020).

*Game abstraction* is commonly used to reduce the size of a game tree prior to solving Billings et al. (2003); Gilpin and Sandholm (2006); Brown and Sandholm (2015b); Čermák, Božansky, and Lisý (2017). Practical abstraction techniques were fundamental to achieving superhuman performance in no-limit Texas hold’em poker in the *Libratus* (Brown and Sandholm 2017b) and *Pluribus* (Brown and Sandholm 2019b) agents. However, these techniques do not have exploitability guarantees. There has been recent work on abstraction algorithms with exploitability guarantees for specific settings (Sandholm and Singh 2012; Basilico and Gatti 2011) and for general extensive-form games (e.g., (Kroer and Sandholm 2014, 2018)), but these are not scalable for large games such as no-limit Texas hold’em, and the guarantees depend on the difference between the abstracted game and the real game being known.

We introduce an approach that can provide exploitability guarantees (even zero exploitability) in black-box games without ever exploring the entire game tree. We introduce a notion of certificate that is often much smaller than the full game. We show that a certificate can be verified in time linear in the size of the certificate, without expanding the remainder of the game tree. For zero-sum games, we give an algorithm that computes an optimal certificate given the current set of explored nodes using any zero-sum game solver as a subroutine. Leveraging prior results, we show that perfect-information (Knuth and Moore 1975) and normal-form (Lipton, Markakis, and Mehta 2003) games have short certificates. We prove that extensive-form games do not always have such, but under a certain informational assumption they do. We also show that it is NP-hard to approximate to within a logarithmic factor the smallest certificate of a game, even in the zero-sum setting, and give an exponential lower bound for the time complexity of solving

86 a black-box game as a function of the size of its smallest  
87 certificate. Despite these hardness results, we give a game-  
88 solving algorithm that expands nodes incrementally until a  
89 certificate is found. It often terminates while only exploring  
90 a small fraction of the tree, and works even when the game  
91 tree is infinite and payoffs may be unbounded. Our experi-  
92 ments show that large and even infinite games can be solved  
93 exactly while expanding only a small fraction of the game  
94 tree.

## 95 2 Preliminaries

96 We study *extensive-form games*, hereafter simply *games*. An  
97 extensive-form game consists of the following:

- 98 (1) a set of players  $\mathcal{P}$ , usually identified with positive inte-  
99 gers  $1, 2, \dots, n$ . *Nature*, a.k.a. *chance*, will be referred to  
100 as player 0. For a given player  $i$ , we will often use  $-i$  to  
101 denote all players except  $i$  and nature.
- 102 (2) a finite tree  $H$  of *histories*, rooted at some *initial state*  
103  $\emptyset \in H$ . The set of leaves, or *terminal states*, in  $H$  will be  
104 denoted  $Z$ . The edges connecting any node  $h \in H$  to its  
105 children are labeled with *actions*.
- 106 (3) a map  $P : H \rightarrow \mathcal{P} \cup \{0\}$ , where  $P(h)$  is the player who  
107 acts at node  $h$  (possibly nature).
- 108 (4) for each player  $i$ , a *utility function*  $u_i : Z \rightarrow \mathbb{R}$ . 5) for  
109 each player  $i$ , a partition of player  $i$ 's decision points, i.e.,  
110  $P^{-1}(i)$ , into *information sets*. In each information set  $I$ ,  
111 every pair of nodes  $h, h' \in I$  must have the same set of  
112 actions. 6) for each node  $h$  at which nature acts, a distribu-  
113 tion  $\sigma_0(\cdot|h)$  over the actions available to nature at node  
114  $h$ .

115 We will use  $(G, u)$ , or simply  $G$  when the utility func-  
116 tion is clear, to denote a game.  $G$  contains the tree and in-  
117 formation set structure, and  $u = (u_1, \dots, u_n)$  is the profile  
118 of utility functions. For any history  $h \in H$  and any player  
119  $i \in \mathcal{P}$ , the *sequence* of player  $i$  at node  $h$  is the sequence of  
120 information sets observed and actions taken by player  $i$  on  
121 the path from the root node to  $h$ . In this paper, all games are  
122 assumed to have perfect recall.

123 A *behavior strategy* (hereafter simply *strategy*)  $\sigma_i$  for  
124 player  $i$  is, for each information set  $I \in J_i$  at which player  
125  $i$  acts, a distribution  $\sigma_i(\cdot|I)$  over the actions available at that  
126 info set. When an agent reaches information set  $I$ , it chooses  
127 action  $a$  with probability  $\sigma_i(a|I)$ .

128 A collection  $\sigma = (\sigma_1, \dots, \sigma_n)$  of behavior strategies, one  
129 for each player  $i \in \mathcal{P}$ , is a *strategy profile*. The *reach prob-*  
130 *ability*  $\sigma_i(h)$  is the probability that node  $h$  will be reached,  
131 assuming that player  $i$  plays according to strategy  $\sigma_i$ , and  
132 all other players (including nature) always choose actions  
133 leading to  $h$  when possible. Analogously, we define  $\sigma(h) =$   
134  $\prod_{i \in \mathcal{P} \cup \{0\}} \sigma_i(h)$  to be the probability that  $h$  is reached under  
135 strategy profile  $\sigma$ . This definition naturally extends to sets of  
136 nodes or to sequences by summing the reach probabilities of  
137 all relevant nodes. A strategy profile induces a distribution  
138 over the terminal nodes of the game. The *value* of a strategy  
139 profile  $\sigma$  for player  $i$  is  $u_i(\sigma) := \mathbb{E}_{z \sim \sigma} u_i(z)$ .

140 The *best response value*  $u_i^*(\sigma_{-i})$  for player  $i$  against an  
141 opponent strategy  $\sigma_{-i}$  is the largest achievable value; i.e. in

142 a two-player game,  $u_i^*(\sigma_{-i}) = \max_{\sigma_i} u_i(\sigma_i, \sigma_{-i})$ . A strat-  
143 egy  $\sigma_i$  is an  $\varepsilon$ -*best response* to opponent strategy  $\sigma_{-i}$  if  
144  $u_i(\sigma_i, \sigma_{-i}) \geq u_i^*(\sigma_{-i}) - \varepsilon$ .

145 A strategy profile  $\sigma$  is an  $\varepsilon$ -*Nash equilibrium* (NE) if all  
146 players are playing  $\varepsilon$ -best responses. *Best responses* and  
147 *Nash equilibria* are respectively 0-best responses and 0-  
148 Nash equilibria.

## 149 3 $\varepsilon$ -Nash certificates via pseudogames

150 We are interested in finding small *certificates* of exact and  
151 approximate Nash equilibria. We introduce a construct that  
152 we call a *pseudogame*, which can be used to build small cer-  
153 tificates of equilibria.

154 **Definition 3.1.** A *pseudogame*  $\tilde{G} = (\tilde{G}, \alpha, \beta)$  is a game in  
155 which some terminal nodes do not have specified utility but  
156 rather have only lower and upper bounds on utilities. For-  
157 mally, for each player  $i$ , instead of the standard utility func-  
158 tion  $u_i : Z \rightarrow \mathbb{R}$ , there are lower and upper bound functions  
159  $\alpha_i : Z \rightarrow \mathbb{R}$  and  $\beta_i : Z \rightarrow \mathbb{R}$  indicating lower and up-  
160 per bounds respectively on the utility of a node. We demand  
161  $\alpha_i(z) \leq \beta_i(z)$  for every  $i$  and  $z$ . We call a node *pseudoter-*  
162 *terminal* if  $\alpha_i(z) < \beta_i(z)$  for some  $i$ , and use *terminal node*  
163 to refer to any leaf in a pseudogame.

164 **Definition 3.2.** An  $\varepsilon$ -Nash equilibrium of a pseudogame  
165  $(\tilde{G}, \alpha, \beta)$  is a strategy profile  $\sigma$  for which, for every player  
166  $i$ , we have  $\beta_i^*(\sigma_{-i}) - \alpha_i(\sigma) \leq \varepsilon$ .

167 **Definition 3.3.** A pseudogame  $(\tilde{G}, \alpha, \beta)$  is a *trunk* of a game  
168  $(G, u)$  if:

- 169 (1)  $\tilde{G}$  can be created by collapsing some internal nodes of  $G$   
170 into terminal nodes (and removing them from information  
171 sets they are contained in), and
- 172 (2) if  $h$  is a pseudoterminal node of  $\tilde{G}$ , and  $z$  is a terminal  
173 node of  $G$  that is a descendant of  $h$ , then  $\alpha_i(h) \leq u_i(z) \leq$   
174  $\beta_i(h)$  for every  $i$ . That is, the bounds  $\alpha$  and  $\beta$  are correct.

175 It is possible for information sets of a game  $G$  to be par-  
176 tially or totally removed in a trunk game.

177 **Definition 3.4.** An  $\varepsilon$ -*certificate* for a game  $G$  is a pair  
178  $(\tilde{G}, \sigma)$ , where  $\tilde{G}$  is a trunk of  $G$  and  $\sigma$  is an  $\varepsilon$ -Nash equi-  
179 librium of  $\tilde{G}$ .

180 Importantly, the definition of a certificate is independent  
181 of the original game  $G$ ; that is, given  $(\tilde{G}, \sigma^*)$ ,  $\varepsilon$  can be com-  
182 puted without knowing the remainder of the game tree of  
183  $G$ : by computing the best response for each player in their  
184 optimistic game, it can be done in time linear in the size of  
185  $\tilde{G}$ .

186 The proposition below shows that our definition of certifi-  
187 cate is reasonable. Proofs are in the appendix.

188 **Proposition 3.5.** Let  $(\tilde{G}, \sigma)$  be an  $\varepsilon$ -certificate for game  $G$ .  
189 Then any strategy profile in  $G$  created by playing according  
190 to  $\sigma$  in any information set appearing in  $\tilde{G}$  and arbitrarily  
191 at information sets not appearing in  $\tilde{G}$  is an  $\varepsilon$ -NE in  $G$ .

## 4 Do small certificates exist?

In this section, we study when games have small  $\varepsilon$ -certificates. Our general goal will be to find certificates of size  $O(N^c \text{poly}(1/\varepsilon))$  for some universal constant  $c < 1$ , where  $N$  is the number of nodes. If a game has a small certificate, there is hope of finding such a certificate quickly, and thus being able to find and verify an (approximate or exact) Nash equilibrium while exploring only a small part of the game. We start by giving a connection between sparse equilibria and small certificates, which we will use later in this section.

**Proposition 4.1** (Sparse equilibria imply small certificates).

Let  $\sigma$  be an  $\varepsilon$ -NE of a game  $G$ , and let  $\tilde{G}$  be the smallest trunk of game  $G$  containing every node  $h$  for which  $\sigma_{-i}(h) > 0$  for any player  $i$ . Then  $(\tilde{G}, \sigma)$  is an  $\varepsilon$ -certificate of  $G$ .

### 4.1 Perfect-information zero-sum games have small certificates, via alpha-beta search

In two-player perfect-information zero-sum games, under certain assumptions, small certificates exist. Specifically, assume that

- (1) there is no randomness (no nature nodes),
- (2) all nodes have uniform branching factor  $b = O(1)$ ,
- (3) moves alternate; i.e., a player-1 decision node is always followed by a player-2 decision node, and
- (4) the tree has uniform depth  $d$ .

In this case, the game has  $N = b^d$  terminal nodes. Alpha-beta search with an optimal heuristic will search only  $O(b^{d/2}) = O(\sqrt{N})$  tree nodes before arriving at a provably optimal strategy (Knuth and Moore 1975). Thus, the portion of the game tree consisting of nodes touched by alpha-beta search contains  $O(\sqrt{N})$  nodes, and constitutes a 0-certificate.

### 4.2 Normal-form games have small certificates, via sparse equilibria

A *normal-form game* is a game in which each player has only a single information set. A two-player normal-form game with  $a_1$  player-1 moves and  $a_2$  player-2 moves (hence  $N = a_1 a_2$  terminal nodes) can thus be expressed as a pair of utility matrices  $A, B \in \mathbb{R}^{a_1 \times a_2}$ . In two-player normal-form games, for every  $\varepsilon$ , there is an  $\varepsilon$ -NE in which each player  $i$  randomizes over  $O(\log(a_{-i})/\varepsilon^2)$  pure strategies Lipton, Markakis, and Mehta (2003). Let  $\sigma^*$  be such an  $\varepsilon$ -Nash equilibrium, and let  $S_i \subseteq [a_i]$  be the support of  $\sigma_i$ .

Consider the following extensive-form pseudogame: First, P1 chooses her strategy  $s_1 \in [a_1]$ . Then, P2 decides whether or not she should play a node from  $S_2$ . If P2 decides not to play from  $S_2$ , and P1 has not played an action in  $S_1$ , the pseudogame terminates immediately in a pseudoterminal node with trivial payoff bounds, i.e.,  $(-\infty, \infty)$ . Otherwise, P2 chooses some strategy  $s_2 \in S_2$  to play, and the proper payoffs are given out. This pseudogame has  $O(a_1|S_2| + a_2|S_1|)$  terminal nodes, and by Proposition 4.1,

the profile  $\sigma^*$  is an  $\varepsilon$ -NE in it. Thus, when  $a_1 = \Theta(a_2)$ , an  $a_1 \times a_2$  normal-form game has an  $\varepsilon$ -certificate of size  $O(\sqrt{N} \log(N)/\varepsilon^2)$ .

Unlike in the case of perfect-information zero-sum games, normal-form games in general do not have small exact certificates: an exact certificate must necessarily include all strategies played in some equilibrium, and there are normal-form games for which the only equilibria are fully mixed.

### 4.3 Extensive-form games with low information have small certificates

This can be generalized to extensive-form games where players do not learn too much information.

**Theorem 4.2.** Let  $G$  be a two-player game with  $N$  nodes and bounded payoffs, and let  $D$  be the maximum number of terminal sequences in the support of any pure strategy for either player. Then  $G$  has an  $\varepsilon$ -Nash equilibrium in which both players mix among  $O((D^2/2\varepsilon^2) \log N)$  pure strategies.

Intuitively,  $D$  is a measure of how much information the players have in the game. A player who learns no information whatsoever throughout the game will have  $D = 1$ , so this proposition matches the sparseness result (Lipton, Markakis, and Mehta 2003) in the normal-form case. On the other hand, a player with perfect information may have  $D = \Omega(\sqrt{N})$  or even larger, in which case this proposition is vacuous.

Under the assumptions of Section 4.1 except perfect information, any given pure strategy is supported on  $O(\sqrt{N})$  nodes. Thus, by Proposition 4.1, we have the following result which implies the existence of small certificates when  $D = O(N^c)$  for  $c < 1/4$ :

**Corollary 4.3.** Under the assumptions of Theorem 4.2 and Section 4.1 except perfect information,  $G$  has an  $\varepsilon$ -certificate of size  $O(\sqrt{N}(D^2/\varepsilon^2) \log N)$ .

As in the case of normal-form games, in general, exact certificates may need to include the whole game tree. However, in some cases, we can do better. For example, games with a natural *public game tree*<sup>1</sup> (Johanson et al. 2011) often have sparse equilibrium strategies (Schmid, Moravcik, and Hladik 2014) and thus small certificates by Proposition 4.1. We will also show later with empirical experiments that many practical games have small exact certificates.

### 4.4 Small certificates do not always exist in extensive-form games

In light of the above results, one might hope that there are sparse approximate equilibria in extensive-form games, which would allow small certificates in such games:

**Question 4.4** (Existence of small  $\varepsilon$ -certificates). Let  $G$  be a two-player zero-sum game with  $N$  nodes. Suppose that  $G$  satisfies the assumptions in Section 4.1. Let  $\varepsilon > 0$ . Is there always an  $\varepsilon$ -certificate with  $O(N^c \text{poly}(1/\varepsilon))$  tree nodes, for some universal constant  $c < 1$ ?

<sup>1</sup>Informally, the public game tree is the game tree visible to an observer with no knowledge of the players' private information.

297 It would be nice if this had a positive answer, since that  
 298 would interpolate between the cases of normal form and  
 299 perfect information, which, as discussed above, both have  
 300  $\tilde{O}(\sqrt{N}/\varepsilon^2)$ -sized certificates. We show that, unfortunately,  
 301 the answer is negative. As a counterexample, consider play-  
 302 ing  $T$  rounds of matching pennies. After each round, P2  
 303 learns what P1 played, but P1 does not learn what P2 played.  
 304 Each round is worth  $1/T$  points, so the maximum score is 1.  
 305 The game tree has uniform depth  $2T$  and uniform branching  
 306 factor 2, for a total of  $N := 2^{2T}$  terminal nodes.

307 **Theorem 4.5.** *Any  $\varepsilon$ -certificate of this game must have at*  
 308 *least  $\Omega(N^{1-O(\varepsilon)})$  nodes.*

309 It does not help to add the assumption that the game  
 310 is win-loss: any zero-sum game can be made win-loss by  
 311 adding normal-form gadget games to the terminal nodes  
 312 which force the players to mix.

## 313 5 Black-box setting

314 For the remainder of this paper, we will assume that we are  
 315 not given access to the full game tree. Instead, we are only  
 316 given black-box access to the game, in the form of a function  
 317 that, given a node  $h$  (in the form of a history of actions),  
 318 gives us:

- 319 (1) upper and lower bounds on the value of any terminal de-  
 320 scendant of  $h$ ,
- 321 (2) if  $h$  is nonterminal, the player to act at that node, and a list  
 322 of legal actions; and
- 323 (3) if the player to act at  $h$  is nature, a single sampled action  
 324 from nature’s action distribution.

325 The game may possibly be very large, or even infinite, but  
 326 we will assume that every node has some terminal descen-  
 327 dant (so that (1) is well-defined), and that the game has a  
 328 finite 0-certificate. The bounds given by (1) may be infinite,  
 329 either because the oracle does not give optimal bounds, or  
 330 because the game is infinite and the payoffs along a branch  
 331 may be unbounded.

332 The first challenge is approximating the true nature distri-  
 333 butions via samples. We thus give a result regarding the sam-  
 334 ple complexity of doing this for a given pseudogame with  
 335 bounded payoffs<sup>2</sup>.

336 **Theorem 5.1** (Sample complexity of approximating a  
 337 game). *Let  $G$  be a game with  $N$  nodes and bounded pay-*  
 338 *offs, and suppose that the true nature distributions are un-*  
 339 *known but have been approximated by sampling at ev-*  
 340 *ery nature node. Let  $\hat{\sigma}_0$  be the approximated nature strat-*  
 341 *egy resulting from this sampling. Fix a player  $i$ . Let  $\hat{u}_i(\sigma)$*   
 342 *denote the expected utility of player  $i$  when the players*  
 343 *play strategy  $\sigma$  and nature plays  $\hat{\sigma}_0$ . Let  $D$  be the max-*  
 344 *imum support size over terminal nodes of any pure strat-*  
 345 *egy profile in the perfect-information refinement of  $G$ . Sup-*  
 346 *pose that, for every nature node  $h$  is sampled at least*  
 347  *$\hat{\sigma}_0(h)(D^2/2\varepsilon^2) \log(2N/\delta)$  times. Then, with probability  $1 -$*   
 348  *$\delta$ , for any strategy profile  $\sigma$ , we have  $|u_i(\sigma) - \hat{u}_i(\sigma)| \leq \varepsilon$ .*

<sup>2</sup>In the unbounded payoff case, the task is hopeless, since it is  
 always possible for there to be a branch of infinite expectation that  
 is reached so rarely that it has never been sampled.

349 Here,  $D$  is some measure of how much randomness there  
 350 is in  $G$ . For example, if  $G$  has no nature nodes,  $D = 1$ . If  $G$   
 351 has no player nodes,  $D = N$ .

352 **Corollary 5.2.** *Let  $\tilde{G}$  be a pseudogame, and consider ap-*  
 353 *proximating nature’s strategy in  $\tilde{G}$  to precision  $\varepsilon$  as per The-*  
 354 *orem 5.1. Let  $\sigma$  be an  $\varepsilon'$ -equilibrium of the approximated*  
 355 *version of  $\tilde{G}$ . Then  $\sigma$  is also an  $(\varepsilon' + 2\varepsilon)$ -equilibrium of  $\tilde{G}$*   
 356 *with probability at least  $1 - 2\delta|\mathcal{P}|$ .*

357 In the above results, the (pseudo)game and sample size  
 358 at each nature node  $h$  are both held fixed; the probability is  
 359 only over the random samples themselves. Thus, if running  
 360 an algorithm that incrementally expands nodes in a pseu-  
 361 dogame, the samples should in principle be re-drawn every  
 362 time  $\tilde{G}$  changes. The factor of  $2|\mathcal{P}|$  is not bothersome since  
 363  $|\mathcal{P}| \leq N$  surely, so this incurs at most a constant factor in  
 364 the sample complexity. Importantly, the sample complexity  
 365 depends only on the size and structure of the pseudogame  
 366  $\tilde{G}$ , not on whatever full game  $G$  that  $\tilde{G}$  may be a trunk of.

367 In the rest of the paper, both for simplicity and to allow  
 368 discussion of the case of unbounded payoffs, we will not  
 369 deal with sampling. Instead, we will assume that the exact  
 370 nature action distribution is given by the black-box oracle  
 371 when a nature node is reached.

## 372 6 The zero-sum case

373 Our results so far have been valid for  $n$ -player general-sum  
 374 games unless otherwise stated. In this section we focus on  
 375 two-player zero-sum games, where one can hope<sup>3</sup> to perhaps  
 376 efficiently *find* small certificates. A two-player game is *zero-*  
 377 *sum* if  $u_1 = -u_2$ . In this case, we refer to a single utility  
 378 function  $u$ ; it is understood that player 2’s utility function is  
 379  $-u$ . In zero-sum games, all Nash equilibria have the same  
 380 expected value; this is called the *value of the game*, and we  
 381 denote it by  $u^*$ . The *exploitability* of an opponent strategy  
 382  $\sigma_{-i}$  for player  $i$  is then  $|u^*(\sigma_{-i}) - u^*|$ .

### 383 6.1 Certificates in zero-sum games

384 In the zero-sum case, we use a slightly different notion of  
 385  $\varepsilon$ -equilibrium of a pseudogame, which will make the subse-  
 386 quent results more precise.

387 **Definition 6.1.** A two-player pseudogame  $(\tilde{G}, \alpha, \beta)$  is zero-  
 388 sum if  $\alpha_2 = -\beta_1$  and  $\beta_2 = -\alpha_1$ .

389 As alluded to above, in this situation, we will drop the  
 390 subscripts, and write  $\alpha$  and  $\beta$  to mean  $\alpha_1$  and  $\beta_1$ . In partic-  
 391 ular,  $(\tilde{G}, \alpha)$  and  $(\tilde{G}, \beta)$  are zero-sum games.

392 **Definition 6.2.** An  $\varepsilon$ -Nash equilibrium of a two-player zero-  
 393 sum pseudogame  $(\tilde{G}, \alpha, \beta)$  is a strategy profile  $(x^*, y^*)$  for  
 394 which  $\beta^*(y^*) - \alpha^*(x^*) \leq \varepsilon$ .

395 In this sense,  $\varepsilon$  is the sum of the exploitabilities of both  
 396 players’ strategies. These are related to Definition 3.2 as fol-  
 397 lows:

<sup>3</sup>In the general-sum setting, finding an approximate Nash equi-  
 librium is PPAD-complete, even for two players (Rubinfeld 2016),  
 so we do not hope to devise certificate-finding algorithms for that  
 case.

398 **Proposition 6.3.** Any  $\varepsilon$ -NE in the sense of Definition 6.2 is  
399 an  $\varepsilon$ -NE in the sense of Definition 3.2.

400 **Proposition 6.4.** Any  $\varepsilon$ -NE in the sense of Definition 3.2 is  
401 a  $2\varepsilon$ -NE in the sense of Definition 6.2.

402 Let  $(\tilde{G}, \alpha, \beta)$  be a pseudogame. Let  $(x_*, y_*)$  be a Nash  
403 equilibrium of the game  $(\tilde{G}, \alpha)$ , and  $(x^*, y^*)$  be a Nash equi-  
404 librium of  $(\tilde{G}, \beta)$ . We will call the pair of strategies  $(x_*, y^*)$   
405 a *pessimistic equilibrium* of  $(\tilde{G}, \alpha, \beta)$  since both players are  
406 playing as if their utilities are as bad as possible. Similarly,  
407 we will call  $(x^*, y_*)$  an *optimistic profile*<sup>4</sup>.

408 By definition, the pessimistic equilibrium is an  $\varepsilon$ -NE of  
409  $(\tilde{G}, \alpha, \beta)$ , where  $\varepsilon = \beta^* - \alpha^*$ . This gives us an algorithm  
410 for finding the best certificate from a given trunk, that runs  
411 in time polynomial in the size of the trunk: to get a strat-  
412 egy for P1 (the maximizer player), solve the game  $(\tilde{G}, \alpha)$ ,  
413 and to get a strategy for P2, solve  $(\tilde{G}, \beta)$ . Since the zero-  
414 sum game solver is used strictly as a subroutine, any solver  
415 of choice may be used: for example, a linear program (LP)  
416 solver with the sequence-form LP (Koller, Megiddo, and von  
417 Stengel 1994; Zhang and Sandholm 2020), modern variants  
418 of CFR (Brown and Sandholm 2019a, 2017a; Brown, Kroer,  
419 and Sandholm 2017; Brown and Sandholm 2015a), or first-  
420 order methods (Hoda et al. 2010; Kroer et al. 2020). If the  
421 solver only finds an  $\varepsilon'$ -equilibrium of the game it is solving,  
422 the result is a certificate for  $(\varepsilon + 2\varepsilon')$ -equilibrium.

## 423 6.2 Lower bounds

424 Since solving zero-sum games can be done efficiently, there  
425 is some hope that small certificates can also be found effi-  
426 ciently. Another goal may be to find a certificate efficiently,  
427 say, in time polynomial in the size of the smallest certificate  
428 of a given game. Unfortunately, these are both impossible:

429 **Theorem 6.5** (Hardness of approximating the smallest cer-  
430 tificate). *Assuming  $P \neq NP$ , there is no  $\text{poly}(N, 1/\varepsilon)$ -time*  
431 *algorithm that, given the game tree of a zero-sum game with*  
432  *$N$  nodes, outputs the smallest  $\varepsilon$ -certificate of the game to*  
433 *better than a  $\Theta(\log N)$  factor of approximation.*

434 **Theorem 6.6.** *There is no algorithm for zero-sum game*  
435 *solving in the black-box setting, even assuming bounded*  
436 *branching factor, with runtime subexponential in the size of*  
437 *the smallest certificate.*

438 These hardness results have slightly different flavors and  
439 consequences. The hardness in Theorem 6.5 comes from  
440 the imperfect information: in the perfect-information set-  
441 ting, the task can be done with a variant of alpha-beta search  
442 in linear time. Further, in practice, we usually do not care  
443 about finding the *smallest* certificate, as long as we can effi-  
444 ciently find one of reasonable size. The hardness in Theo-  
445 rem 6.6 is more fundamental: it comes from the fact that we  
446 cannot assume access to any reasonable heuristic of where  
447 to explore; thus, we may explore the optimal path of play  
448 last in the worst case, resulting in a large certificate.

<sup>4</sup>The pessimistic equilibrium is an equilibrium of the pseu-  
dogame. The optimistic profile may not be, hence the difference  
in naming.

## 6.3 An algorithm for solving black-box games

450 Despite the difficulties presented by Theorems 6.5 and 6.6,  
451 we present an algorithm for finding a certificate in a zero-  
452 sum game in the black-box setting, with nontrivial provable  
453 guarantees. For now, we will assume that the game  $\tilde{G}$  has  
454 bounded payoffs; later we will relax this assumption.

---

**Algorithm 6.7** Finding a certificate in a two-player zero-  
sum game

---

- 1: start with a pseudogame  $(\tilde{G}, \alpha, \beta)$  that has only a root node.
  - 2: **loop**
  - 3: solve  $(\tilde{G}, \alpha)$  and  $(\tilde{G}, \beta)$  with an LP solver to obtain equilibria  $(x_*, y_*)$  and  $(x^*, y^*)$ .
  - 4: expand all pseudoterminal nodes of  $\tilde{G}$  that appear in the support of  $(x^*, y_*)$ .
  - 5: (if there are none, stop and output  $\tilde{G}$  and the pessimistic equilibrium  $(x_*, y^*)$ .)
- 

455 We use LP for the game solves in Line 3, for three reasons. First, LP<sup>5</sup> results in an exact solution (at least up to numerical tolerances), which is desirable because the support of the solution is relevant to Line 4; iterative solvers such as CFR typically return fully mixed solutions. Second, only a small number of changes are made to the LP with each node expanded, so LP algorithms that can be warm started, such as primal or dual simplex, can be efficient in practice. Third, it will allow us to adapt this algorithm to the case of unbounded payoffs, which we will see later; again, CFR cannot do that.

456 From the discussion in Section 6.1, we know that this algorithm will always output an 0-certificate. If we want an  $\varepsilon$ -certificate for  $\varepsilon > 0$ , we can also simply terminate the algorithm when  $\beta^* - \alpha^* < \varepsilon$ . We now prove an important fact about Algorithm 6.7.

457 **Theorem 6.8.** *A pseudogame has a 0-Nash equilibrium if and only if it has an optimistic profile with no pseudoterminal node in its support.*

458 The “only if” direction guarantees that Line 4 does not terminate the algorithm unless a 0-certificate has been found. The “if” direction guarantees a weak form of “this algorithm will not waste work”: modulo the uniqueness of the optimistic profile<sup>6</sup>, the algorithm stops exactly when it has found a 0-certificate. This is not trivial: other protocols such as “expand all pseudoterminal nodes appearing in the support of at least one player in the pessimistic equilibrium” fail to satisfy the “if” direction.

459 The algorithm has no runtime bound as a function of the size of the smallest certificate of  $G$ , even assuming bounded

---

<sup>5</sup>using either an exact method such as simplex, or an interior-point method such as barrier with crossover

<sup>6</sup>When the optimistic profile is not unique, the algorithm may waste work: for example, there may be one equilibrium which has support over pseudoterminal nodes and one which does not, the algorithm may pick the former and continue expanding nodes, making an unnecessarily big (but still correct) certificate.

485 branching factor: indeed, if  $G$  is infinite, it is even possible  
 486 for the algorithm to run indefinitely, even when a finite-sized  
 487 certificate exists. One way to fix this without losing more  
 488 than a constant factor in efficiency is to, in addition to Line 4,  
 489 also always expand the shallowest strictly pseudoterminal  
 490 node of  $\tilde{G}$  at each iteration. This way, a certificate with  $d$   
 491 nodes has depth at most  $d$ , and thus will be generated after  
 492 at most after  $O(b^d)$  expansions (where  $b$  is a bound on the  
 493 branching factor of the game), matching the lower bound of  
 494 Theorem 6.6.

## 495 6.4 Handling unbounded payoffs

496 In infinite games with unbounded payoffs, it is possible for  
 497 the games  $(\tilde{G}, \alpha)$  and  $(\tilde{G}, \beta)$  to have infinite-magnitude util-  
 498 ity on some nodes. For example,  $(\tilde{G}, \beta)$  may have payoff  
 499  $+\infty$  on some nodes (but not  $-\infty$ ). We now show how to  
 500 adapt Algorithm 6.7 for such situations. Assume WLOG  
 501 that we are solving  $(\tilde{G}, \beta)$ ; i.e. it is possible for payoffs  
 502 to be  $+\infty$  but not  $-\infty$  (for  $(\tilde{G}, \alpha)$ , swap the players).  
 503 Call a P2-sequence *bad* if its support (over terminal nodes)  
 504 contains a node of utility  $+\infty$ . Assume that it is possi-  
 505 ble for P2 to avoid all bad sequences; otherwise, the game  
 506 has value  $+\infty$ . Consider the sequence-form bilinear saddle-  
 507 point problem (Koller, Megiddo, and von Stengel 1994)  
 508 for  $(\tilde{G}, \beta)$  (Equation (6.9)) and its equivalent LP (Equa-  
 509 tion (6.10)):

$$\max_{x \geq 0} \min_{y \geq 0} x^T A y \quad \text{s.t. } Bx = b, Cy = c, x, y \geq 0 \quad (6.9)$$

$$\max_{x \geq 0, z} c^T z \quad \text{s.t. } Bx = b, C^T z \leq A^T x. \quad (6.10)$$

510 Here  $A$  is the payoff matrix, which may contain infinite en-  
 511 tries. Then, the main idea is to remove any constraint cor-  
 512 responding to bad P2-sequences, and solve the resulting LP  
 513 (which now by construction contains no infinite entries and  
 514 is thus well formed), for a Nash equilibrium solution  $x$ . The  
 515 problem is that  $x$  may not be a true Nash equilibrium of  
 516  $(\tilde{G}, \beta)$ , since it is possible for P1 to end up avoiding nodes  
 517 of utility  $+\infty$ , which could allow P2 to best respond by ac-  
 518 tually playing toward a bad sequence.

519 Let  $V^*(s)$  denote the value that P2 receives by playing  
 520 a best response to  $x$  starting at a P2 infoset or sequence  $s$ .  
 521 Let  $V(s)$  denote the same, except while forcing P2 to avoid  
 522 bad sequences. Obviously,  $V^* \leq V$ . Consider the following  
 523 recursive algorithm, which we run on every P2-root infoset  
 524  $I$ :

---

**Algorithm 6.11** CORRECT( $I$ ): Correcting a strategy in the  
 case of infinite reward

---

- 1: **for** each action  $a$  available to P2 **do**
  - 2:   **if**  $V^*(Ia) < V(I)$  **then**
  - 3:     **for** every P1-sequence  $i$  such that  $A_{i,Ia} = +\infty$  **do**  
        $x_i \leftarrow x_i + \text{infinitesimal}$ <sup>7</sup>
  - 4:     **for** every P2-infoset  $I'$  whose parent sequence is  
        $Ia$  **do** CORRECT( $I'$ )
- 

<sup>7</sup>This can be easily formalized by perturbing by  $\varepsilon$ , then taking

525 Call a pair of strategies a *corrected optimistic profile* if it  
 526 is the result of applying this procedure to both parts of an  
 527 optimistic profile. We can now make the following strength-  
 528 ening of Theorem 6.8:

**Theorem 6.12.** *A pseudogame with possibly unbounded  
 529 payoffs has a 0-Nash equilibrium if and only if it has a cor-  
 530 rected optimistic profile with no pseudoterminal node in its  
 531 support.* 532

533 Thus, to run Algorithm 6.7 in games with unbounded pay-  
 534 offs, it suffices to apply the correction algorithm to the opti-  
 535 mistic profile found in Line 3 before expanding nodes.

## 536 7 Experiments

537 We conducted experiments using the algorithm in Section 6  
 538 on the following common zero-sum benchmark games.

- (1) A zero-sum variant of the **search game** Bořanský and  
 539 Čermák (2015). 540
- (2)  **$k$ -rank Goofspiel**. It is played as follows. At time  $t$  (for  
 541  $t = 1, \dots, k$ ), players place bids for a prize of value  $t$ . The  
 542 possible bids are the integers  $1, \dots, k$ , and each player  
 543 must bid each integer exactly once. The player with the  
 544 higher bid wins the prize; if the bids are equal, the prize  
 545 is split equally. The winner of each round is made public  
 546 after each round, but the bids are not. The goal of each  
 547 player is to maximize the sum of the values of her prizes  
 548 won. In the *perfect-information* (PI) variant, P2 knows  
 549 P1's bid while bidding, and bids are made public after  
 550 each round. This creates a perfect-information game in  
 551 which P2 has a large advantage, and in which we expect  
 552 a certificate of size  $O(\sqrt{N})$ . In the *random* variant, the  
 553 order of the prizes is randomized. 554
- (3)  **$k$ -rank limit Leduc poker**. It is a small variant of limit  
 555 poker, played with one hole card and one community card,  
 556 and a deck with  $k$  ranks. The players are only allowed  
 557 to raise by a fixed amount, but can do so an unlimited  
 558 number of times. Thus, the possible payoffs in the game,  
 559 and the length of the game, are both unbounded. 560

561 We computed 0-certificates in all cases. For the LP solver,  
 562 we used Gurobi v9.0.0 (Gurobi Optimization, LLC 2019).  
 563 Results of experiments can be found in Table 1. In many  
 564 games, we found 0-certificates of size substantially smaller  
 565 than the number of nodes in the game, and the certificate size  
 566 as a fraction of the game size decreases as the game grows.

567 The results in Goofspiel align with the theoretical pre-  
 568 dictions: perfect-information games have very small certifi-  
 569 cates (basically  $\sqrt{N}$  nodes). In light of Proposition 4.1, it  
 570 also makes sense that certificates are smaller (relative to the  
 571 size of the game) when there is no randomness: randomness  
 572 simply increases the number of nodes in the game tree repre-  
 573 sented by any given pure strategy, so an equilibrium with the  
 574 same sparsity for the players now leads to a larger certificate.

$\varepsilon$  sufficiently small. The strategy need not actually ever be con-  
 575 structed, so there is no need to formally discuss how small  $\varepsilon$   
 576 needs to be; if coding this algorithm, we can simply store the indices of  
 577 infinitesimal entries.

Table 1: Experimental results. The *minimal certificate* is a certificate after removing all unnecessary nodes per Proposition 4.1. Percentages are relative to game size. Leduc variants have infinite size; for them, “game size” reported is for the trunk with the number of consecutive raises restricted to 12.

game	size of game		size of certificate				size of minimal certificate			
	nodes	infosets	nodes		infosets		nodes		infosets	
search game	234,705	11,890	13,682	5.8%	532	4.5%	5,526	2.4%	379	3.2%
4-rank PI Goofspiel	2,229	1,653	275	12.3%	110	6.7%	141	6.3%	54	3.3%
5-rank PI Goofspiel	55,731	41,331	2,593	4.7%	957	2.3%	763	1.4%	288	0.7%
6-rank PI Goofspiel	2,006,323	1,487,923	21,948	1.1%	7,584	0.5%	4,438	0.2%	1,677	0.1%
4-rank Goofspiel	2,229	738	614	27.5%	117	15.9%	294	13.2%	58	7.9%
5-rank Goofspiel	55,731	9,948	11,415	20.5%	2,160	21.7%	8,518	15.3%	1,792	18.0%
6-rank Goofspiel	2,006,323	166,002	266,756	13.3%	15,776	9.5%	171,343	8.5%	12,135	7.3%
3-rank random Goofspiel	1,066	426	309	29.0%	92	21.6%	214	20.1%	65	15.3%
4-rank random Goofspiel	68,245	17,432	16,416	24.1%	3,270	18.8%	11,992	17.6%	2,335	13.4%
5-rank random Goofspiel	8,530,656	1,175,330	1,854,858	21.7%	241,985	20.6%	1,388,172	16.3%	185,946	15.8%
5-rank limit Leduc	197,736	13,920	26,306	13.3%	2,406	17.3%	12,923	6.5%	1,242	8.9%
9-rank limit Leduc	1,181,512	44,928	137,662	11.7%	6,811	15.2%	51,533	4.4%	2,891	6.4%
13-rank limit Leduc	3,578,472	93,600	337,312	9.4%	12,171	13.0%	105,769	3.0%	4,449	4.8%

In Leduc poker, no node involving more than 12 consecutive raises was ever expanded in any size of game while searching for a certificate. This suggests that it is never optimal for either player to play past this point, despite the fact that continuing to raise could in principle lead to an unbounded payoff. This phenomenon allows our algorithm to find a finite-sized 0-certificate, thus completely solving the game in a reasonably efficient manner, even though it has infinite size.

## 8 Conclusions and future research

We presented a notion of certificate for general extensive-form games that allows verification of exact and approximate Nash equilibria without expanding the whole game tree. We showed that small equilibria exist in some restricted classes of extensive-form game, but not all. We presented algorithms for both verifying a certificate and computing the optimal certificate given the currently-explored trunk of a game. Our experiments showed that many large or even infinite games have small certificates, allowing us to find equilibria while exploring a vanishingly small portion of the game.

This paper opens many directions for future research:

- (1) Develop further the ideas of Section 5 for the case of unknown nature distributions. For example, what is the best way to balance sampling, game tree exploration, and equilibrium finding?
- (2) Seek algorithms for finding certificates that give stronger guarantees of optimality than Theorem 6.12, especially in the case of infinite games with unbounded utilities.
- (3) Seek algorithms with stronger guarantees than that implied by Proposition 4.1 for verifying the Nash gap of a given strategy profile; for example, is it possible to eas-

ily construct the smallest trunk for which a given  $\sigma$  is an  $\epsilon$ -equilibrium?

## Broader Impacts

The techniques have broad applicability. Furthermore, the paper opens up additional important research directions.

Improving the strategic capabilities of people and companies will typically (but not always) improve systemwide good as the players will be able to better reach win-win solutions. In zero-sum games this is not the case because the size of the “cake” is constant, so there are winners and losers. In both the general case and the zero-sum case, AI tools like the ones in this paper can help elevate less educated and less experienced players up to the same level as expert players, thereby making the distribution of value more fair.

A potential downside is that *if* the technology were only available to the privileged, that could increase unfairness.

## Acknowledgements

This material is based on work supported by the National Science Foundation under grants IIS-1718457, IIS-1617590, IIS-1901403, and CCF-1733556, and the ARO under awards W911NF1710082 and W911NF2010081.

## References

- 628  
629 Basilico, N.; and Gatti, N. 2011. Automated Abstractions for  
630 Patrolling Security Games. In *AAAI Conference on Artificial*  
631 *Intelligence (AAAI)*.
- 632 Berner, C.; Brockman, G.; Chan, B.; Cheung, V.; Debiak, P.;  
633 Dennison, C.; Farhi, D.; Fischer, Q.; Hashme, S.; Hesse, C.;  
634 et al. 2019. Dota 2 with Large Scale Deep Reinforcement  
635 Learning. *arXiv preprint arXiv:1912.06680* .
- 636 Billings, D.; Burch, N.; Davidson, A.; Holte, R.; Schaeffer,  
637 J.; Schauenberg, T.; and Szafron, D. 2003. Approximating  
638 Game-Theoretic Optimal Strategies for Full-scale Poker. In  
639 *Proceedings of the International Joint Conference on Artificial*  
640 *Intelligence (IJCAI)*.
- 641 Bošanský, B.; and Čermák, J. 2015. Sequence-form al-  
642 gorithm for computing Stackelberg equilibria in extensive-  
643 form games. In *Twenty-Ninth AAAI Conference on Artificial*  
644 *Intelligence*.
- 645 Bowling, M.; Burch, N.; Johanson, M.; and Tammelin, O.  
646 2015. Heads-up Limit Hold'em Poker is Solved. *Science*  
647 347(6218).
- 648 Brown, N.; Kroer, C.; and Sandholm, T. 2017. Dynamic  
649 Thresholding and Pruning for Regret Minimization. In *AAAI*  
650 *Conference on Artificial Intelligence (AAAI)*.
- 651 Brown, N.; and Sandholm, T. 2015a. Regret-Based Prun-  
652 ing in Extensive-Form Games. In *Proceedings of the An-*  
653 *ual Conference on Neural Information Processing Systems*  
654 *(NIPS)*.
- 655 Brown, N.; and Sandholm, T. 2015b. Simultaneous Abstrac-  
656 tion and Equilibrium Finding in Games. In *Proceedings of*  
657 *the International Joint Conference on Artificial Intelligence*  
658 *(IJCAI)*.
- 659 Brown, N.; and Sandholm, T. 2017a. Reduced Space and  
660 Faster Convergence in Imperfect-Information Games via  
661 Pruning. In *International Conference on Machine Learning*  
662 *(ICML)*.
- 663 Brown, N.; and Sandholm, T. 2017b. Superhuman AI for  
664 heads-up no-limit poker: Libratus beats top professionals.  
665 *Science* eaa01733.
- 666 Brown, N.; and Sandholm, T. 2019a. Solving imperfect-  
667 information games via discounted regret minimization. In  
668 *AAAI Conference on Artificial Intelligence (AAAI)*.
- 669 Brown, N.; and Sandholm, T. 2019b. Superhuman AI for  
670 multiplayer poker. *Science* 365(6456): 885–890.
- 671 Čermák, J.; Bošanský, B.; and Lisý, V. 2017. An algorithm  
672 for constructing and solving imperfect recall abstractions of  
673 large extensive-form games. In *Proceedings of the Inter-*  
674 *national Joint Conference on Artificial Intelligence (IJCAI)*,  
675 936–942.
- 676 Gilpin, A.; and Sandholm, T. 2006. A Competitive Texas  
677 Hold'em Poker Player via Automated Abstraction and Real-  
678 Time Equilibrium Computation. In *Proceedings of the Na-*  
679 *tional Conference on Artificial Intelligence (AAAI)*, 1007–  
680 1013.
- Gurobi Optimization, LLC. 2019. Gurobi Optimizer Refer- 681  
ence Manual. 682
- Heinrich, J.; and Silver, D. 2016. Deep reinforcement learn- 683  
ing from self-play in imperfect-information games. *arXiv*  
684 *preprint arXiv:1603.01121* . 685
- Hoda, S.; Gilpin, A.; Peña, J.; and Sandholm, T. 2010. 686  
Smoothing Techniques for Computing Nash Equilibria of  
687 Sequential Games. *Mathematics of Operations Research*  
688 35(2). 689
- Johanson, M.; Waugh, K.; Bowling, M.; and Zinkevich, M. 690  
2011. Accelerating Best Response Calculation in Large Ex-  
691 tensive Games. In *Proceedings of the International Joint*  
692 *Conference on Artificial Intelligence (IJCAI)*. 693
- Knuth, D. E.; and Moore, R. W. 1975. An analysis of alpha- 694  
beta pruning. *Artificial Intelligence* 6(4): 293–326. 695
- Koller, D.; Megiddo, N.; and von Stengel, B. 1994. Fast 696  
algorithms for finding randomized strategies in game trees.  
697 In *Proceedings of the 26<sup>th</sup> ACM Symposium on Theory of*  
698 *Computing (STOC)*. 699
- Kroer, C.; and Sandholm, T. 2014. Extensive-Form Game 700  
Abstraction With Bounds. In *Proceedings of the ACM Con-*  
701 *ference on Economics and Computation (EC)*. 702
- Kroer, C.; and Sandholm, T. 2018. A Unified Framework 703  
for Extensive-Form Game Abstraction with Bounds. In *Pro-*  
704 *ceedings of the Annual Conference on Neural Information*  
705 *Processing Systems (NIPS)*. 706
- Kroer, C.; Waugh, K.; Kılınç-Karzan, F.; and Sandholm, T. 707  
2020. Faster algorithms for extensive-form game solving via  
708 improved smoothing functions. *Mathematical Programming*  
709 . 710
- Lanctot, M.; Waugh, K.; Zinkevich, M.; and Bowling, M. 711  
2009. Monte Carlo Sampling for Regret Minimization in  
712 Extensive Games. In *Proceedings of the Annual Conference*  
713 *on Neural Information Processing Systems (NIPS)*. 714
- Lanctot, M.; Zambaldi, V.; Gruslys, A.; Lazaridou, A.; 715  
Tuyls, K.; Pérolat, J.; Silver, D.; and Graepel, T. 2017. A uni-  
716 fied game-theoretic approach to multiagent reinforcement  
717 learning. In *Proceedings of the Annual Conference on Neu-*  
718 *ral Information Processing Systems (NIPS)*, 4190–4203. 719
- Lipton, R.; Markakis, E.; and Mehta, A. 2003. Playing 720  
Large Games Using Simple Strategies. In *Proceedings of*  
721 *the ACM Conference on Electronic Commerce (ACM-EC)*,  
722 36–41. San Diego, CA: ACM. 723
- Marchesi, A.; Trovò, F.; and Gatti, N. 2020. Learn- 724  
ing Probably Approximately Correct Maximin Strategies in  
725 Simulation-Based Games with Infinite Strategy Spaces. In  
726 *Autonomous Agents and Multi-Agent Systems*, 834–842. 727
- Moravčík, M.; Schmid, M.; Burch, N.; Lisý, V.; Morrill, D.; 728  
Bard, N.; Davis, T.; Waugh, K.; Johanson, M.; and Bowling,  
729 M. 2017. DeepStack: Expert-level artificial intelligence in  
730 heads-up no-limit poker. *Science* . 731
- Raz, R.; and Safra, S. 1997. A Sub-Constant Error- 732  
Probability Low-Degree Test, and a Sub-Constant Error-  
733 Probability PCP Characterization of NP. In *Proceedings*  
734



735 of the *Annual Symposium on Theory of Computing (STOC)*,  
736 475–484.

737 Rubinstein, A. 2016. Settling the complexity of computing  
738 approximate two-player Nash equilibria. In *Proceedings of*  
739 *the Annual Symposium on Foundations of Computer Science*  
740 *(FOCS)*, 258–265.

741 Sandholm, T.; and Singh, S. 2012. Lossy stochastic game  
742 abstraction with bounds. In *Proceedings of the ACM Con-*  
743 *ference on Electronic Commerce (EC)*.

744 Schmid, M.; Moravcik, M.; and Hladik, M. 2014. Bound-  
745 ing the support size in extensive form games with imperfect  
746 information. In *AAAI Conference on Artificial Intelligence*  
747 *(AAAI)*, 784–790.

748 Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.;  
749 Van Den Driessche, G.; Schrittwieser, J.; Antonoglou, I.;  
750 Panneershelvam, V.; Lanctot, M.; et al. 2016. Mastering the  
751 game of Go with deep neural networks and tree search. *Nature*  
752 *529(7587)*: 484.

753 Srinivasan, S.; Lanctot, M.; Zambaldi, V.; Pérolat, J.; Tuyls,  
754 K.; Munos, R.; and Bowling, M. 2018. Actor-critic pol-  
755 icy optimization in partially observable multiagent environ-  
756 ments. In *Proceedings of the Annual Conference on Neural*  
757 *Information Processing Systems (NIPS)*, 3422–3435.

758 Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.;  
759 Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds,  
760 T.; Georgiev, P.; et al. 2019. Grandmaster level in Star-  
761 Craft II using multi-agent reinforcement learning. *Nature*  
762 *575(7782)*: 350–354.

763 Wellman, M. 2006. Methods for Empirical Game-Theoretic  
764 Analysis (Extended Abstract). In *Proceedings of the Na-*  
765 *tional Conference on Artificial Intelligence (AAAI)*, 1552–  
766 1555.

767 Zhang, B. H.; and Sandholm, T. 2020. Sparsified Linear  
768 Programming for Zero-Sum Equilibrium Finding. In *Inter-*  
769 *national Conference on Machine Learning (ICML)*.

770 Zhou, Y.; Li, J.; and Zhu, J. 2020. Posterior sampling  
771 for multi-agent reinforcement learning: solving extensive  
772 games with imperfect information. In *International Con-*  
773 *ference on Learning Representations*.

775 **A.1 Proposition 3.5**

$$u_i^*(\sigma_{-i}) - u_i(\sigma) \leq \beta_i^*(\sigma_{-i}) - \alpha_i(\sigma) \leq \varepsilon. \quad \square$$

776 **A.2 Proposition 4.1**

777 By definition, it is impossible to reach any pseudoterminal node of  $\tilde{G}$  by changing only a single player's strategy. Thus, for  
778 any player  $i$ , we have  $\beta_i^*(\sigma_{-i}) - \alpha(\sigma) \leq u_i^*(\sigma_{-i}) - u(\sigma) \leq \varepsilon$ . (the first inequality may not be an equality, because the best  
779 response  $\beta_i^*(\sigma_{-i})$  is taken in the pseudogame, and  $u_i^*$  is taken in the full game, where there is more flexibility.  $\square$ )

780 **A.3 Theorem 4.5**

781 **Lemma A.1.** *In every  $\varepsilon$ -NE of  $G$ , the entropy of P1's strategy is at least  $T(1 - 2\varepsilon)$  bits.*

782 *Proof.* Let  $\sigma_1$  be any P1 strategy in  $\varepsilon$ -equilibrium, and let  $H_T$  be the entropy over terminal nodes when P1 plays  $\sigma_1$  and P2  
783 plays uniformly at random. Let  $U_T$  be the number of rounds that P2 loses if she best responds to P1. Since  $\sigma_1$  is an  $\varepsilon$ -NE  
784 strategy, we have  $U_T \geq T(1/2 - \varepsilon)$ . We will show that  $H_T \geq 2U_T + T$ , which will complete the proof.

785 Proceed by induction on  $T$ . For  $T = 1$ , the claim follows from the inequality  $h(p) \geq 2 \min(p, 1 - p)$ , which is true for all  
786  $p \in [0, 1]$ , where  $h$  is the binary entropy function.

787 In the inductive case, suppose that, at the top information set, P1 plays strategy  $x = [p, q]$  (i.e. heads with probability  $p$ , and  
788 tails with probability  $q$ ). Let  $H' \in \mathbb{R}^{2 \times 2}$  be the matrix whose  $ij$ -entry is the conditional entropy over terminal nodes after P1  
789 plays  $i$  and P2 plays  $j$  in the root information set. Similarly, let  $U'$  be the matrix of conditional remaining expected number of  
790 rounds lost, not including this round, for player 2. Note that the utility matrix of the overall game, assuming that P2 plays  
791 correctly in later rounds, is  $A := U' + I$ . By IH,  $H' \geq 2U' + T - 1$  element-wise. Further, P2's move in this information set  
792 does not affect the future of the game, since P1 does not learn P2's move, and P2's move does not otherwise affect her future  
793 optimal decisions. That is,  $U'y$  is the same for all (normalized)  $y$ . Let  $y$  be the uniform random strategy for player 1, and  $y^*$   
794 be a best response for player 1. Then we have:

$$\begin{aligned} H &= 1 + h(p) + x^T H' y \\ &\geq T + h(p) + 2x^T U' y \\ &= T + h(p) + 2x^T U' y^* \\ &= T + h(p) + 2x^T A y^* - 2x^T y^* \\ &= T + h(p) + 2x^T A y^* - 2 \min(p, 1 - p) \end{aligned}$$

795 and we are once again done by the inequality  $h(p) \geq 2 \min(p, 1 - p)$ .  $\square$

796 The restriction on P2's strategy is necessary: indeed, since P1 has only  $2^T$  pure strategies, there are sparse  $\varepsilon$ -NE strategies  
797 for P2 supported on only  $O(T/\varepsilon^2)$  pure strategies.

798 Somewhat surprisingly, this proposition becomes false if P1 learns what P2 played in each round. Indeed, the P1 strategy  
799 "play heads if your number of losses minus number of wins is  $\varepsilon T$ , and uniformly at random otherwise" is (for large  $T$ ) an  
800  $\varepsilon$ -equilibrium with basically  $T$  bits of entropy, since if P2 plays uniformly at random, with very good probability their score  
801 delta will never exceed  $\varepsilon T$ . However, despite having low entropy, this strategy has a very large support over terminal nodes.

802 **Corollary A.2.** *In every  $\varepsilon$ -NE of this game, for every  $t \geq T/2$ , the first  $t$  rounds of P1's strategy have at least  $t(1 - 4\varepsilon)$  bits of  
803 entropy.*

804 **Corollary A.3.** *Let  $\varepsilon \leq 1/16$ . In every  $\varepsilon$ -NE of this game, for every  $t \geq T/2$ , P1's strategy assigns probability at least  $2^{-t}$  to  
805 at least half of her pure strategies at round  $t$ .*

806 *Proof.* Let  $Z$  be a random variable for P1's selected strategy, and  $E$  be the event that  $Z$  is among the half least likely pure  
807 strategies to be picked.

$$H(Z) = H(Z, E) = \Pr[E]H(Z|E) + \Pr[\neg E]H(Z|\neg E) + H(E) \leq \frac{2^t p t}{2} + \frac{t}{2}$$

808 where  $H$  is the entropy. We know from above that  $H(Z) \geq t(1 - 4\varepsilon)$ , so the claim follows by solving for  $p$ .  $\square$

809 We now prove Theorem 4.5. The proof acts like a partial converse to Proposition 4.1 for this game. Let  $((\tilde{G}, \alpha, \beta), \sigma)$  be an  
810  $\varepsilon$ -certificate, and let  $Z'$  be the set of terminal nodes in  $\tilde{G}$ . Let  $u$  be the assignment of utilities induced by P2 playing uniform

random at every decision point outside  $\tilde{G}$  (it does not matter at this point how P1 plays). Let  $\sigma'_i$  be the uniform random strategy for player  $i$ . Then: 811

$$\beta_2(\sigma_1, \sigma'_2) \leq \beta_2^*(\sigma_1) \leq u_2(\sigma) + \varepsilon \leq u_2(\sigma'_1, \sigma_2) + 2\varepsilon = u_2(\sigma_1, \sigma'_2) + 2\varepsilon. \quad (\text{A.4}) \quad 812$$

For simplicity of notation, for any terminal node  $z$  of  $\tilde{G}$ , let  $r(z)$  be the number of rounds remaining in the game. Then note that  $\beta(z) - u(z) = r(z)/2T$  for every  $z$ . Now suppose for contradiction that  $\tilde{G}$  has fewer than  $n := 2^{2T(1-16\varepsilon)-2}$  terminal nodes. Consider the level of the game tree after both players have made  $t := (1 - 16\varepsilon)T$  moves; in other words, the level at which  $r(z) = 16\varepsilon T$ . This level has  $4n$  nodes, so certainly  $\tilde{G}$  must contain at most  $1/4$  of the nodes at this level. Let  $S$  be a set of half of the nodes of  $G$  at level  $t$  to which P1 assigns probability at least  $2^{-t}$ . Then  $\tilde{G}$  contains at most half the nodes in  $S$ . Now observe that 813

$$\begin{aligned} \beta_2(\sigma_1, \sigma'_2) - u_2(\sigma_1, \sigma_2^*) &= \frac{1}{2T} \mathbb{E}_z r(z) \\ &\geq \frac{1}{2T} \sum_{z \in S \setminus \tilde{G}} \sigma_1(z) \sigma_2^*(z) r(z) \\ &\geq \frac{1}{2T} \frac{1}{2} 2^{2t} 2^{-t} r(z) = 4\varepsilon \end{aligned} \quad 814$$

which contradicts (A.4). 815

#### A.4 Theorem 4.2 816

We first introduce some terminology that will be useful in this section. The *realization plan* corresponding to a strategy  $\sigma_i$  is the vector of reach probabilities  $\sigma_i(s)$  for each *sequence*  $s$  for player  $i$ . The constraints on valid realization plans are linear, and the payoff of a two-player zero-sum game can be expressed as a bilinear form  $x^T A y$ , where  $x$  and  $y$  are the realization plan vectors for the two players, and  $A$  is a payoff matrix depending only on the terminal node values (Koller, Megiddo, and von Stengel 1994). This bilinear program is known as the *sequence form* of a game. 817

**Lemma A.5.** *Let  $x$  be any P1 strategy. Let  $\hat{x}$  be a strategy profile defined by mixing uniformly at random over a multiset of  $k$  independent sampled pure strategies from  $x$ , where* 818

$$k \geq \frac{D^2}{2\varepsilon^2} \log \frac{2N}{\delta}. \quad 819$$

and  $D$  is the maximum support size over terminal sequences of any P2 pure strategy. Then with probability  $1 - \delta$ , for any strategy profile  $y$ , we have  $|u_2(\hat{x}, y) - u_2(x, y)| \leq \varepsilon$ . 820

*Proof.* We follow basically the same idea as the proof in Lipton, Markakis, and Mehta (2003). Let  $A$  be the P2 sequence-form payoff matrix, restricted to those rows and columns corresponding to terminal sequences. By Hoeffding, we have 821

$$\Pr \left[ |(A\hat{x})_i - (Ax)_i| \geq \frac{\varepsilon}{D} \right] \leq 2e^{-2k\varepsilon^2/D^2} \leq \frac{\delta}{N}$$

by picking  $k$  as above. Taking a union bound over the at most  $N$  sequences for P2, we have  $\|A\hat{x} - Ax\|_\infty \leq \varepsilon/D$  with probability  $1 - \delta$ . Now select an  $x'$  for which this is true. Then by Hölder's inequality, for any pure realization plan  $y$ , we have 822

$$|y^T A\hat{x} - y^T Ax| \leq \|y\|_1 \|A\hat{x} - Ax\|_\infty \leq \varepsilon. \quad 823$$

where the last inequality follows because  $\|y\|_1 \leq D$ . Now since  $|y^T A\hat{x} - y^T Ax|$  is convex in  $y$ , and the pure realization plans are the vertices of the polytope of all realization plans, we are done. 824

Theorem 4.2 now follows by applying the lemma to an equilibrium strategy  $x$  with any  $\delta < 1$ . 825

#### A.5 Theorem 5.1 826

Sampling this number of samples at each nature node  $h$  is at least as good as sampling  $(D^2/2\varepsilon^2) \log(2N/\delta)$  pure nature strategies. The proposition now follows by applying Lemma A.5 to the game in which the game tree is the same as  $G$ , P1 is nature, P2 controls every actual player in  $G$  (and thus has perfect information), and the P2 utility function is  $u$ . 827

#### A.6 Corollary 5.2 828

By a union bound over the  $|\mathcal{P}|$  players and the two utility functions  $\alpha_i$  and  $\beta_i$  for each player, we have that with probability at least  $1 - 2\delta|P|$ , for every  $i$  and every deviation  $\sigma'_i$ ,  $|\hat{\alpha}_i(\sigma'_i, \sigma_{-i}) - \alpha_i(\sigma'_i, \sigma_{-i})| \leq \varepsilon$  and  $|\hat{\beta}_i(\sigma'_i, \sigma_{-i}) - \beta_i(\sigma'_i, \sigma_{-i})| \leq \varepsilon$ . 829

Let  $\hat{\alpha}_i(\sigma)$  and  $\hat{\beta}_i(\sigma)$  for a given strategy  $\sigma$  be the utilities of  $\sigma$  under the approximated version of  $\tilde{G}$ . Let  $\hat{\sigma}_i^*$  be a best response for player  $i$  in the approximated version of  $\tilde{G}$ , and let  $\sigma_i^*$  be a best response in  $\tilde{G}$  itself. Then we have: 830

$$\beta_i^*(\sigma_{-i}) \leq \hat{\beta}_i(\sigma_i^*, \sigma_{-i}) + \varepsilon \leq \hat{\beta}_i^*(\sigma_{-i}) + \varepsilon \leq \hat{\alpha}_i(\sigma) + \varepsilon + \varepsilon' \leq \alpha_i(\sigma) + 2\varepsilon + \varepsilon' \quad 831$$

for every player  $i$ . 832

847 **A.7 Proposition 6.3**

848 Let  $(x, y)$  be an  $\varepsilon$ -NE in the sense of Definition 6.2. Then

$$\beta^*(y) - \alpha(x, y) \leq \beta^*(y) - \alpha^*(x) \leq \varepsilon \quad \text{and} \quad \beta(x, y) - \alpha^*(x) \leq \beta^*(y) - \alpha^*(x) \leq \varepsilon. \quad \square$$

849 **A.8 Proposition 6.4**

850 Let  $(x, y)$  be an  $\varepsilon$ -NE in the sense of Definition 3.2. Then

$$\beta^*(y) - \alpha^*(x) \leq \beta^*(y) - \alpha(x, y) + \beta(x, y) - \alpha^*(x) \leq 2\varepsilon. \quad \square$$

851 **A.9 Theorem 6.5**

852 We reduce from the SET-COVER problem, which is known to be NP-hard to better than a  $\Theta(\log n)$  factor (Raz and Safra  
853 1997). In SET-COVER, we are given a universe  $U = \{1, \dots, n\}$  and a collection of  $m$  sets  $\mathcal{S} = \{S_1, \dots, S_m\}$  whose union is  
854  $U$ , and our task is to find the smallest subset of  $\mathcal{S}$  whose union is still  $U$ .

855 Consider the following game: P2 starts by choosing to either *play* or *leave*. If P2 leaves, then the game immediately termi-  
856 nates, and P1 gets value  $1/2m$ . If P2 chooses to play, then P1 chooses an index  $i = 1, \dots, m$ . Then, P1 is given  $m$  consecutive  
857 opportunities to leave the game (and immediately lose), should they choose. (The sole purpose of this is to inflate the size of  
858 the certificate.) After this, P2, without knowing the  $i$ , chooses an element  $u \in U$ . P1 gets value 1 if  $u \in S_i$ , and 0 otherwise.

859 This game has  $\text{poly}(m, n)$  nodes, and its value (for P1) is exactly  $1/2m$ , since P1 can force P2 to leave by playing uniformly  
860 at random (and not choosing to lose). We now claim that, for  $\varepsilon < 1/2m$ , finding an  $\varepsilon$ -certificate of size  $\Theta((m+n)k)$  is  
861 equivalent to finding a set cover of size  $k$ , which completes the proof.

862 If  $\mathcal{R} \subseteq \mathcal{S}$  is a set cover of size  $k$ , then consider the trunk created by expanding exactly those P2 decision nodes where P1  
863 has played some set  $S_i \in \mathcal{R}$ . This creates a trunk of size  $\Theta((m+n)k)$ . Even pessimistically, P1 can gain value  $1/k \geq 1/m$  by  
864 randomizing uniformly over  $\mathcal{R}$  in this trunk; thus, P2 is forced to leave, and this is a 0-certificate.

865 Conversely, suppose we had an  $\varepsilon$ -certificate, for  $\varepsilon < 1/2m$ , constructed from some tree  $\tilde{G}$ . Let  $\mathcal{R}$  be the collection of sets  
866  $S_i \in \mathcal{S}$  for which P2's decision node after P1 plays  $S_i$  has been expanded, and let  $k = |\mathcal{R}|$ . Then the trunk has size at least  
867  $\Omega((m+n)k)$ . If  $\mathcal{R}$  is not a set cover, then there is some  $u \in U$  outside the union of sets in  $\mathcal{R}$ . If P1 plays  $u$ , then she gains  
868 optimistic value 0. Thus, since  $\varepsilon < 1/2m$ ,  $\mathcal{R}$  must be a set cover.  $\square$

869 **A.10 Theorem 6.6**

870 Consider the family of two-player games in which there is a target string  $x \in \{0, 1\}^n$ , and play proceeds as follows: Player 1  
871 chooses, bit-by-bit, a string  $y \in \{0, 1\}^n$ . If  $x = y$ , then Player 1 wins; otherwise, Player 2 chooses whether to win or lose. The  
872 smallest certificate in this game has size  $\Theta(n)$ , and consists of the path of play to  $y$ . However, there is no algorithm, randomized  
873 or deterministic, that will find the correct node  $y$  without first expanding  $\Omega(2^n)$  other nodes.  $\square$

874 **A.11 Theorem 6.8**

875 ( $\Leftarrow$ ) Suppose  $\tilde{G}$  has no 0-certificate. Let  $(x^*, y_*)$  be an optimistic profile. Then

$$\alpha(x^*, y_*) \leq \alpha^*(y_*) < \beta^*(x^*) \leq \beta(x^*, y_*).$$

876 where the middle inequality is strict since  $\tilde{G}$  has no 0-certificate. But then  $\alpha(x^*, y_*) \neq \beta(x^*, y_*)$ ; i.e., there is some uncertainty  
877 as to the value of the strategy profile  $(x^*, y_*)$ ; i.e., there is a nonzero probability that a pseudoterminal node is reached.

878 ( $\Rightarrow$ ) Now suppose  $\tilde{G}$  has a 0-certificate, and call it  $(x_*, y^*)$ . Clearly,  $(x_*, y^*)$  cannot contain in its support any pseudoterminal  
879 node. We claim that  $(x_*, y^*)$  is also an optimistic profile of  $\tilde{G}$ , which completes the proof. Indeed, we have

$$\alpha^*(x_*) \leq \beta^*(x_*) \leq \beta^*(y^*) \quad \text{and} \quad \alpha^*(x_*) \leq \alpha^*(x^*) \leq \beta^*(y^*)$$

880 But all of these must actually be equalities, since  $\alpha^*(x_*) = \beta^*(y^*)$  for a 0-certificate. Thus,  $x_*$  is a Nash equilibrium strategy  
881 in  $(\tilde{G}, \beta)$ , and  $y^*$  is a Nash equilibrium strategy in  $(\tilde{G}, \alpha)$ , which is what we needed to show.  $\square$

882 **A.12 Theorem 6.12**

883 ( $\Leftarrow$ ) The correction algorithm adds infinitesimal amounts to sequences such that P2 is then forced to never play to any bad  
884 sequence that could be used to achieve value better than  $V(I)$ . Thus, corrected equilibrium is actually an  $\varepsilon$ -equilibrium for  
885 infinitesimal  $\varepsilon$ , and the proof of Appendix A.11 applies verbatim.

886 ( $\Rightarrow$ ) A pessimistic strategy will never be corrected, since a pessimistic player never has a terminal node of utility  $+\infty$ . Thus,  
887 again, the proof of Appendix A.11 applies verbatim.  $\square$