# Hardness of Independent Learning and Sparse Equilibrium Computation in Markov Games

Dylan J. Foster <sup>1</sup> Noah Golowich <sup>2</sup> Sham M. Kakade <sup>3</sup>

#### **Abstract**

We consider the problem of decentralized multiagent reinforcement learning in Markov games. A key question is whether there are algorithms that, when run independently by all agents, lead to noregret for each player, analogous to celebrated results for no-regret learning in normal-form games. While recent work has shown that such algorithms exist for restricted settings (e.g., when regret is defined with respect to deviations to Markov policies), the question of whether independent no-regret learning can be achieved in the standard Markov game framework was open. We provide a decisive negative resolution to this problem, both from a computational and statistical perspective. We show that:

- Under the assumption that PPAD-hard problems cannot be solved in polynomial time, there is no polynomial-time algorithm that attains noregret in general-sum Markov games when executed independently by all players, even when the game is known to the algorithm designer and the number of players is a small constant.
- 2. When the game is unknown, no algorithm, efficient or otherwise, can achieve no-regret without observing exponentially many episodes in the number of players.

These results are proven via lower bounds for a simpler problem we refer to as Sparsecce, in which the goal is to compute a coarse correlated equilibrium that is "sparse" in the sense that it can be represented as a mixture of a small number of product policies.

Proceedings of the 40 <sup>th</sup> International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

#### 1. Introduction

The framework of multi-agent reinforcement learning (MARL), which describes settings in which multiple agents interact in a dynamic environment, has played a key role in recent breakthroughs in artificial intelligence, including the development of agents that approach or surpass human performance in games such as Go (Silver et al., 2016), Poker (Brown & Sandholm, 2018), Stratego (Perolat et al., 2022), and Diplomacy (Kramar'et al., 2022; Bakhtin et al., 2022). MARL also shows promise for real-world multi-agent systems, including autonomous driving (Shalev-Shwartz et al., 2016), and cybersecurity (Malialis & Kudenko, 2015), and economic policy (Zheng et al., 2022). These applications, where reliability is critical, necessitate the development of algorithms that are practical and efficient, yet provide strong formal guarantees and robustness.

Multi-agent reinforcement learning is typically studied using the framework of Markov games (also known as stochastic games) (Shapley, 1953). In a Markov game, agents interact over a finite number of steps: at each step, each agent observes the state of the environment, takes an action, and observes a reward which depends on the current state as well as the other agents' actions. Then the environment transitions to a new state as a function of the current state and the actions taken. An episode consists of a finite number of such steps, and agents interact over the course of multiple episodes, progressively learning new information about their environment. Markov games generalize the well-known model of Markov Decision Processes (MDPs) (Puterman, 1994), which describe the special case in which there is a single agent acting in a dynamic environment, and we wish to find a policy that maximizes its reward. By contrast, for Markov games, we typically aim to find a distribution over agents' policies which constitutes some type of equilibrium.

#### 1.1. Decentralized learning

In this paper, we focus on the problem of decentralized (or, independent) learning in Markov games. In decentralized MARL, each agent in the Markov game behaves independently, optimizing their policy myopically while treating the effects of the other agents as exogenous. Agents observe local information (in particular, their own actions and rewards), but do not observe the actions of the other agents directly. Decentralized learning enjoys a number of desirable properties, including

<sup>&</sup>lt;sup>1</sup>Microsoft Research <sup>2</sup>Massachusetts Institute of Technology <sup>3</sup>Harvard University. Correspondence to: Dylan J. Foster<dylanfoster@microsoft.com>, Noah Golowich<nzg@mit.edu>, Sham M. Kakade<sham@seas.harvard.edu>.

scalability, versatility, and practicality. The central question we consider is whether there exist decentralized learning algorithms which, when employed by all agents in a Markov game, lead them to play near-equilibrium strategies over time.

Decentralized equilibrium computation in MARL is not well understood theoretically, and algorithms with provable guarantees are scarce. To motivate the challenges and most salient issues, it will be helpful to contrast with the simpler problem of decentralized learning in normal-form games, which may be interpreted as Markov games with a single state. Much of the modern work on decentralized learning in normal-form games centers on no-regret learning, where agents select actions independently using online learning algorithms (Cesa-Bianchi & Lugosi, 2006) designed to minimize their regret (that is, the gap between realized payoffs and the payoff of the best fixed action in hindsight). In particular, a foundational result is that if each agent employs a no-regret learning strategy, then the average of the agents' joint action distributions approaches a coarse correlated equilibrium (CCE) for the normal-form game (Cesa-Bianchi & Lugosi, 2006; Hannan, 1957; Blackwell, 1956). CCE is a natural relaxation of the foundational concept of Nash equilibrium, which has the downside of being intractable to compute. On the other hand, there are many efficient algorithms that can achieve vanishing regret in a normal-form game, even when opponents select their actions in an arbitrary, potentially adaptive fashion, and thus converge to a CCE (Vovk, 1990; Littlestone & Warmuth, 1994; Cesa-Bianchi et al., 1997; Hart & Mas-Colell, 2000; Syrgkanis et al., 2015).

This simple connection between no-regret learning and decentralized convergence to equilibria has been influential in game theory, leading to numerous lines of research including fast rates of convergence to equilibria (Syrgkanis et al., 2015; Chen & Peng, 2020; Daskalakis et al., 2021; Anagnostides et al., 2022), price of anarchy bounds for smooth games (Roughgarden, 2015), and lower bounds on query and communication complexity for equilibrium computation (Fearnley et al., 2013; Rubinstein, 2016; Babichenko & Rubinstein, Empirically, no-regret algorithms such as regret matching (Hart & Mas-Colell, 2000) and Hedge (Vovk, 1990; Littlestone & Warmuth, 1994; Cesa-Bianchi et al., 1997) have been used to compute equilibria that can achieve state-of-the-art performance in application domains such as Poker (Brown & Sandholm, 2018) and Diplomacy (Bakhtin et al., 2022). Motivated by these successes, we ask whether an analogous theory can be developed for Markov games. In particular:

Are there efficient algorithms for no-regret learning in Markov games?

Challenges for no-regret learning. In spite of active research effort and many promising pieces of progress (Jin et al., 2021; Song et al., 2022; Mao & Basar, 2021; Daskalakis et al., 2022; Erez et al., 2022), no-regret learning guarantees for Markov games have been elusive. A barrier faced by naive algorithms

is that it is intractable to ensure no-regret against an arbitrary adversary, both computationally (Bai et al., 2020; Abbasi Yadkori et al., 2013) and statistically (Liu et al., 2022; Kwon et al., 2021; Foster et al., 2022). Fortunately, many of the implications of no-regret learning (in particular, convergence to equilibria) do not require the algorithm to have sublinear regret against an arbitrary adversary, but rather only against other agents who are running the same algorithm independently. This observation has been influential in normal-form games, where the line of work on fast rates of convergence to equilibrium (Syrgkanis et al., 2015; Chen & Peng, 2020; Daskalakis et al., 2021; Anagnostides et al., 2022) holds only in this more restrictive setting. This motivates the following relaxation to our central question.

Problem 1.1. Is there an efficient algorithm that, when adopted by all agents in a Markov game and run independently, leads to sublinear regret for each individual agent?

Attempts to address Problem 1.1. Two recent lines of research have made progress toward addressing Problem 1.1 and related questions. In one direction, several recent papers have provided algorithms, including V-learning (Jin et al., 2021; Song et al., 2022; Mao & Basar, 2021) and SPoCMAR (Daskalakis et al., 2022), that do not achieve no-regret, but can nevertheless compute and then sample from a coarse correlated equilibrium in a Markov game in a (mostly) decentralized fashion, with the caveat that they require a shared source of random bits as a mechanism to coordinate. Notably, V-learning depends only mildly on the shared randomness: agents first play policies in a fully independent fashion (i.e., without shared randomness) according to a simple learning algorithm for T episodes, and use shared random bits only once learning finishes as part of a post-processing procedure to extract a CCE policy. A question left open by these works, is whether the sequence of policies played by the V-learning algorithm in the initial independent phase can itself guarantee each agent sublinear regret.

Most closely related to our work, Erez et al. (2022) recently showed that Problem 1.1 can be solved positively for a restricted setting in which regret for each agent is defined as the maximum gain in value they can achieve by deviating to a fixed Markov policy. Markov policies are those whose choice of action depends only on the current state as opposed to the entire history of interaction. This notion of deviation is restrictive because in general, even when the opponent plays a sequence of Markov policies, the best response will be non-Markov. In challenging settings that abound in practice, it is standard to consider non-Markov policies (Leibo et al., 2021; Agapiou et al., 2022), since they often achieve higher value than Markov policies; we provide a simple example in Proposition B.1. Thus, while a regret guarantee with respect to the class of Markov policies (as in (Erez et al., 2022)) is certainly interesting, it may be too weak in general, and it is of great interest to understand whether Problem 1.1 can be answered positively in the general setting.<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>We remark that the V-learning and SPoCMAR algorithms

We refer the reader to Appendix B.2 for further discussion.

#### 1.2. Our contributions

We resolve Problem 1.1 in the negative, from both a computational and statistical perspective.

Computational hardness. We provide two computational lower bounds (Theorems 1.2 and 1.3) which show that under standard complexity-theoretic assumptions, there is no efficient algorithm that runs for a polynomial number of episodes and guarantees each agent non-trivial ("sublinear") regret when used in tandem by all agents. Both results hold even if the Markov game is explicitly known to the algorithm designer; Theorem 1.3 is stronger and more general, but applies only to 3-player games, while Theorem 1.2 applies to 2-player games, but only for agents restricted to playing Markovian policies.

To state our first result, Theorem 1.2, we define a product Markov policy to be a joint policy in which players choose their actions independently according to Markov policies (see Sections 2 and 3 for formal definitions). Note that if all players use independent no-regret algorithms to choose Markov policies at each episode, then their joint play at each round is described by a product Markov policy, since any randomness in each player's policy must be generated independently.

Theorem 1.2 (Informal version of Corollary 3.3). If PPAD P, then there is no polynomial-time algorithm that, given the description of a 2-player Markov game, outputs a sequence of joint product Markov policies which guarantees each agent sublinear regret.

Theorem 1.2 provides a decisive negative resolution to Problem 1.1 under the assumption that PPAD P, <sup>2</sup> which is standard in the theory of computational complexity (Papadimitriou, 1994). <sup>3</sup> Beyond simply ruling out the existence of fully decentralized no-regret algorithms, it rules out existence of centralized algorithms that compute a sequence of product policies for which each agent has sublinear regret, even if such a sequence does not arise naturally as the result of agents independently following some learning algorithm. Salient implications include:

 Theorem 1.2 provides a separation between Markov games and normal-form games, since standard no-regret algorithms for normal-form games i) run in polynomial

mentioned above do learn equilibria that are robust to deviations to non-Markov policies, though they do not address Problem 1.1 since they do not have sublinear regret.

<sup>2</sup>Technically, the class we are denoting by P, namely of total search problems that have a deterministic polynomial-time algorithm, is sometimes denoted by FP, as it is a search problem. We ignore this distinction.

<sup>3</sup> PPAD is the most well-studied complexity class in algorithmic game theory, and is widely believed to not admit polynomial time algorithms. Notably, the problem of computing a Nash equilibrium for normal-form games with two or more players is PPAD-complete (Daskalakis et al., 2009; Chen et al., 2006; Rubinstein, 2018).

time and ii) produce sequences of joint product policies that guarantee each agent sublinear regret. Notably, no-regret learning for normal-form games is efficient whenever the number of agents is polynomial, whereas Theorem 1.2 rules out polynomial-time algorithms for as few as two agents.

A question left open by the work of Jin et al. (2021); Song et al. (2022); Mao & Basar (2021) was whether the sequence of policies played by the V-learning algorithm during its independent learning phase can guarantee each agent sublinear regret. Since V-learning plays product Markov policies during the independent phase and is computationally efficient, Theorem 1.2 implies that these policies do not enjoy sublinear regret (assuming PPAD P).

Our second result, Theorem 1.3, extends the guarantee of Theorem 1.2 to the more general setting in which agents can select arbitrary, potentially non-Markovian policies at each episode. This comes at the cost of only providing hardness for 3-player games as opposed to 2-player games, as well as relying on the slightly stronger complexity-theoretic assumption that PPAD † RP.4

Theorem 1.3 (Informal version of Corollary 4.4). If PPAD † RP, then there is no polynomial-time algorithm that, given the description of a 3-player Markov game, outputs a sequence of joint product general policies (i.e., potentially non-Markov) which guarantees each agent sublinear regret.

Statistical hardness. Theorems 1.2 and 1.3 rely on the widely-believed complexity theoretic assumption that PPAD-complete problems cannot be solved in (randomized) polynomial time. Such a restriction is inherent if we assume that the game is known to the algorithm designer. To avoid complexity-theoretic assumptions, we consider a setting in which the Markov game is unknown to the algorithm designer, and algorithms must learn about the game by executing policies ("querying") and observing the resulting sequences of states, actions, and rewards. Our final result, Theorem 1.4, shows unconditionally that, for m-player Markov games whose parameters are unknown, any algorithm computing a no-regret sequence as in Theorem 1.3 requires a number of queries that is exponential in m.

Theorem 1.4 (Informal version of Theorem 5.2). Given query access to a m-player Markov game, no algorithm that makes fewer than 2  $\,^{\rm pmq}$  queries can output a sequence of joint product policies which guarantees each agent sublinear regret.

Similar to our computational lower bounds, Theorem 1.4 goes far beyond decentralized algorithms, and rules out even centralized algorithms that compute a no-regret sequence by jointly controlling all players. The result provides another

<sup>4</sup>We use RP to denote the class of total search problems for which there exists a polynomial-time randomized algorithm which outputs a solution with probability at least 2{3, and otherwise outputs "fail".

separation between Markov games and normal-form games, since standard no-regret algorithms for normal-form games can achieve sublinear regret using polypmq queries for any m.

The 2

pmq scaling in the lower bound, which does not rule out query-efficient algorithms when m is constant, is to be expected for an unconditional result: If the game has only polynomially many parameters (which is the case for constant m), one can estimate all of the parameters using standard techniques (Jin et al., 2020), then directly find a no-regret sequence.

Proof techniques: the SPARSECCE problem. Our proofs proceed via establishing lower bounds for a computational problem we refer to as SPARSECCE. In the SPARSECCE problem, the aim is to compute a CCE that can be represented as the mixture of a small number of product policies. See Sections 3 and 4 for detailed proof overview.

Organization. Section 2 presents preliminaries, Sections 3 and 4 provide our computational lower bounds, and Section 5 presents our unconditional lower bounds for multi-player games.

Notation. For nPN, we write rns:t1;2;:::;nu. For a finite set T, pT q denotes the space of distributions on T. For an ele-ment tPT, It PpT q denotes the delta distribution that places probability mass 1 on t. We adopt standard big-oh notation, and write f rOpgq to denote that f Opg maxt1;polylogpgquq,

with r
pq and
pq defined analogously.

#### 2. Preliminaries

This section contains preliminaries necessary to present our main results. We first introduce the Markov game framework (Section 2.1), then provide a brief review of normal-form games (Section 2.3), and finally introduce the concepts of coarse correlated equilibria and regret minimization (Section 2.4).

#### 2.1. Markov games

We consider general-sum Markov games in a finite-horizon, episodic framework. For mPN, an m-player Markov game G consists of a tuple G pS;H;pA<sub>i</sub>q<sub>iPrms</sub>;P;pR<sub>i</sub>q<sub>iPrms</sub>;q, where:

- S denotes a finite state space and H P N denotes a finite time horizon. We write S:|S|.
- For i P rms,  $A_i$  denotes a finite action space for agent i. We let  $A: {}^{m-1}A_1^i$  denote the joint action space and  $A_i: {}^{i_1}A_1^i$ . We denote joint actions in bold, e.g., a pa<sub>1</sub>;:::;a<sub>m</sub>q PA. We write  $A_i:|A_i|$  and A:|A|.
- P pP<sub>1</sub>; ::: ; P<sub>H</sub>q is the transition kernel, with each P<sub>h</sub> : S A  $\tilde{N}$  pS q denoting the kernel for step h P rHs. In particular, P<sub>h</sub>ps<sup>1</sup>|s;aq is the probability of transitioning to s<sup>1</sup> from the state s at step h when agents play a.
- For i Prms and h PrHs, R<sub>i;h</sub>: S A Ñ r1{H;1{Hs

is the reward function for agent i: 5 the reward agent i receives in state s at step h if agents play a is R<sub>i;h</sub>ps;aq.6

• PpSq denotes the initial state distribution.

An episode in the Markov game proceeds as follows: the initial state  $s_1$  is drawn from the initial state distribution . Then, for each  $h \neq H$ , given state  $s_h$ , each agent i plays action  $a_{i;h}$  P  $A_i$ , and given the joint action profile  $a_h$   $pa_{1;h}$ ; ...;  $a_{m;h}q$ , each agent i receives reward of  $r_{i;h}$   $R_{i;h}ps_h$ ;  $a_hq$  and the state of the system transitions to  $s_h$  1  $P_hp|s_h$ ;  $a_hq$ . We denote the tuple of agents' rewards at each step h by  $r_h$   $pr_{1;h}$ ;...;  $r_{m;h}q$ , and refer to the resulting sequence  $_H$  : $ps_1$ ;  $a_1$ ;  $r_1q$ ;...;  $ps_H$ ;  $a_H$ ;  $r_Hq$  as a trajectory. For h P rHs, we define the prefix of the trajectory via  $_h$  : $ps_1$ ;  $a_1$ ;  $r_1q$ ;...;  $ps_h$ ;  $a_h$ ;  $r_hq$ .

We use the following notation: for some quantity x (e.g., action, reward, etc.) indexed by agents, i.e.,  $x px_1; ...; x_mq$ , and an agent i P rms, we write  $x_i px_1; ...; x_{i1}; x_{i1}; ...; x_mq$  to denote the tuple consisting of all  $x_{i1}$  for  $i^1i$ .

#### 2.2. Policies and value functions

We now introduce the notion of policies and value functions for Markov games. Policies are mappings from states (or sequences of states) to actions for the agents. We consider several different types of policies, which play a crucial role in distinguishing the types of equilibria that are tractable and those that are intractable to compute efficiently.

Markov policies. A randomized Markov policy for agent i is a sequence  $_{i}$   $p_{i;1}$ ;:::; $_{i;H}$ q, where  $_{i;h}$ : S  $\tilde{N}$   $pA_{i}$ q. We denote the space of randomized Markov policies for agent i by  $^{markov}$ . We write  $^{markov}$ : $^{markov}$   $^{markov}$  to denote the space of product Markov policies, which are joint policies in which each agent i independently follows a policy in  $^{markov}$ . In particular, a policy  $^{markov}$  is specified by a collection  $^{p_1}$ ;:::; $^{h}$ q, where  $^{h}$ : S  $\tilde{N}$   $^{pA_1}$ q  $^{pA_m}$ q. We additionally define  $^{markov}$ :  $^{1}$   $^{markov}$ , and for a policy  $^{p}$   $^{markov}$ , write  $^{i}$  to denote, the collection of mappings  $^{i}$   $^{p_{i;1}}$ ;:::; $^{i;H}$ q, where  $^{i;h}$ : S  $\tilde{N}$   $^{i_1}$  $^{i_1}$  $^{pA_{i_1}}$ q denotes the tuple of all but player i's policies.

When the Markov game G is clear from context, for a policy P  $^{markov}$  we let Prs denote the law of the trajectory when players select actions via  $a_h$  pshq, and let Ers denote the corresponding expectation.

General (non-Markov) policies. In addition to Markov policies, we will consider general history-dependent (or, non-

<sup>5</sup>We assume that rewards lie in r1{H;1{Hs for notational convenience, as this ensures that the cumulative reward for each — episode lies in r1;1s. This assumption is not important to our results.

<sup>6</sup>We restrict our attention to Markov games in which the rewards at each step are a deterministic function of the state and action profile. Since our goal is to prove lower bounds, this is without loss.

Markov) policies, which select actions based on the entire sequence of states and actions observed up the current step. To streamline notation, for i P rms, let  $_{i;h}$  ps $_{1;a;;1};r_{i;1};...;s_h;a_{i;h};r_{i;h}q$  denote the history of agent i's states, actions, and reward up to step h. Let  $H_{i;h}$  pS  $A_i$  r0;1sqh denote the space of all possible histories of agent i up to step h. For i P rms, a randomized general (i.e., non-Markov) policy of agent i is a collection of mappings  $_i$  p $_{i;1};...;j_{i;H}q$  where  $_{i;h}$ :  $H_{i;h1}$  S  $\tilde{N}$  pA $_i$ q is a mapping that takes the history observed by agent i up to step h 1 and the current state and outputs a distribution over actions for agent i.

gen;rnd We denote by the space of randomized general policies of agent i, and further write gen;rnd : gen;rnd gen;rnd to denote the space of product general policies; note that markov € gen; and m In particular, a policy P gen;rnd is specicfied by a collection  $p_{i;h}q_{iPrms;hPrHs}$ , where  $i;h:H_{i;h1}$  S  $\tilde{N}$   $pA_iq$ . When agents play according to a general policy Pgen;rnd, at each step h, each agent, given the current state sh and their history  $_{i;h1}$  P  $H_{i;h1}$ , chooses to play an action  $a_{i;h}$   $_{i;h}p_{i;h1}$ ;  $s_hq$ , independently from all other agents. For a policy P gen;rnd, we let Prs and Ers denote the law and expectation operator for the trajectory when players select actions via a<sub>h</sub> p<sub>h1</sub>;s<sub>h</sub>q, and write i to denote the collection of policies of all agents but i, i.e., i pj;hqhPrHs;jPrmsztiu.

We will also consider distributions over product randomized general policies, namely elements of  $p^{gen;rnd}q$ . We will refer to elements of  $p^{gen;rnd}q$  as distributional policies. To play a distributional policy P P  $p^{gen;rnd}q$ , agents draw a randomized policy P (so that  $P^{gen;rnd}$ ) and then play .

Value functions. For a general policy P  $^{gen;rnd}$ , we define the value function for agent i P rms as  $V_i$ : E  $^{H}_{h1}R_{i;h}ps_h;a_hq$  |  $s_1$ ; this represents the expected reward that agent i receives when each agent chooses their actions via  $a_{i;h}$   $_{h}p_{i;h1};s_hq.For$  a distributional policy P Pp $^{gen;rnd}q$ , we extend this notation by defining  $V^P$ :  $E_P rV s$ .

#### 2.3. Normal-form games

To motivate the solution concepts we consider for Markov games, let us first revisit the notion of normal-form games, which may be interpreted as Markov games with a single state. For m; n P N, an m-player n-action normal-form game G is specified by a tuple of m reward tensors  $M_1; ...; M_m P r0; 1s^{nn}$ , where each tensor is of order m (i.e., has  $n^m$  entries). We will write  $GpM_1; ...; M_mq$ . We as-sume for simplicity that each player has the same number n of actions, and identify each player's action space with rns. Then

<sup>7</sup>When T is not a finite set, we take pT q to be the set of Radon probability measures over T equipped with the Borel -algebra.

an an action profile is specified by a Prns<sup>m</sup>; if each player acts according to a, then the reward for player i Prms is given by pM<sub>i</sub>q<sub>a</sub> Pr0;1s. Our hardness results will use the standard notion of Nash equilibrium in normal-form games. We define the mplayer pn;q-NASH problem to be the problem of computing an approximate Nash equilibrium of a given m-player n-action normal-form game. (See Definition C.2 for a formal definition of -Nash equilibrium.) A celebrated result is that Nash equilibria are PPAD-hard to approximate, i.e., the 2-player pn;n<sup>c</sup>q-NASH problem is PPAD-hard for any constant c i 0 (Daskalakis et al., 2009; Chen et al., 2006). We refer the reader to Section C.2 for further background on these concepts.

#### 2.4. Markov games: Equilibria and no-regret

We now turn our focus back to Markov games, and introduce the main solution concepts we consider, as well as the notion of no-regret. Since computing Nash equilibria is intractable even for normal-form games, much of the work on efficient equilibrium computation has focused on alternative notions of equilibrium, notably coarse correlated equilibria.

For a distributional policy P  $Pp^{gen;rnd}q$  and a randomized policy  ${}^1P_i^{gen;rnd}$  of player i, we let  ${}^1P_iP_p^{gen;rnd}q$  denote the distributional policy which is given by the distribution of  $p^1;_iqP^{gen;rnd}$  for P (and i denotes the marginal of on all players but i). For  $P^{gen;rnd}$ , we write  $P^1:_iqP^{gen;rnd}$  be using  $P^1:_iqP^{gen;rnd}$ . Let us fix a Markov game G, which in particular determines the players' value functions V.

Definition 2.1 (Coarse correlated equilibrium). For j = 0, a distributional policy P P pgen,rndq is defined to be an -coarse correlated equilibrium (CCE) if for each i P rms, it holds that  $\max_{i,pgen,rnd} V^{\frac{1}{p_i}} V^{\frac{p}{i}} / .$ 

Coarse correlated equilibria can be computed efficiently for both normal-form games and Markov games, and are fundamentally connected to the notion of no-regret and independent learning, which we now introduce.

Regret. For a policy P<sup>gen;rnd</sup>, we denote the distributional policy which puts all its mass on by I Pp<sup>gen;rnd</sup>q. Thus  $\frac{1}{t^{1}}$  T  $t^{1}$  I<sub>ptq</sub> P p<sup>gen;rnd</sup>q denotes the distributional policy which randomizes uniformly over the ptq. We define regret as follows.

Definition 2.2 (Regret). Consider a sequence of policies p1q;:::;pTq Pgen;rnd. For iPrms, the regret of agent i with respect to this sequence is defined as:

It is immediate from the above definitions that a sequence of policies page:::;p<sup>rq</sup> P gen;rnd satisfies Reg;::p<sup>p1q</sup>;:::;p<sup>rq</sup> q/

T if and only if the distributional policy :  $\frac{1}{2}$  T  $I_{ptq}$  is an -CCE (stated formally in Fact C.1 in the appendix). T

No-regret learning. A standard approach to decentralized equilibrium computation, which exploits Fact C.1, is to select  $^{p1q}$ ;:::; $^{pTq}$  P  $^{gen;rnd}$  using independent no-regret learning algorithms. A no-regret learning algorithm for player i selects  $^{ptq}$  P  $^{gen;rnd}$  based on the realized trajectories  $^{p1q}$  ::; $^{pTq}$  P  $^{ptq}$ ; H that player i observes over the course of play, but with no knowledge of  $^{ptq}$ , so as to ensure that no-regret is achieved: Reg  $_{i;T}$  p $^{p1q}$ ;:::; $^{pTq}$ q/T. If each player i uses their own, independent no-regret learning algorithm, this approach yields product policies  $^{ptq}$  p $^{tq}$  p $^{tq}$ , and the uniform average of the  $^{ptq}$  yields a CCE as long as all of the players can keep their regret small.

For the special case of normal-form games, there are several efficient algorithms, which—when run independently—ensure that each player's regret after T episodes is bounded above by Op Tq (that is Op1{ Tq), even when the other players' actions are chosen adversarially.

#### 3. Lower bound for Markovian algorithms

In this section we prove Theorem 1.2 (restated formally below as Theorem 3.2), establishing that in two-player Markov games, there is no computationally efficient algorithm that computes a sequence play:::;p<sup>Tq</sup> of product Markov policies so that each player has small regret under this sequence. This section serves as a warm-up for our results in Section 4, which remove the assumption that play:::;p<sup>Tq</sup> are Markovian.

#### 3.1. Sparse MarkovCCE and computational model

As discussed in the introduction, our lower bounds for no-regret learning are a consequence of lower bounds for the SPAR-SECCE problem. In what follows, we formalize this problem (specifically, the Markovian variant, which we refer to as SPARSEMARKOVCCE), as well as our computational model.

Description length for Markov games (constant m). Given a Markov game G, we let pGq denote the maximum number of bits needed to describe any of the rewards  $R_{i;h}ps$ ; aq or transition probabilities  $P_hps^1|s$ ; aq in binary. We define  $|G|: maxtS; max_{iPrms}A_i; H; pGqu$ . The interpretation of |G| depends on the number of players m: If m is a constant (as will be the case in the current section and Section 4), then

|G| should be interpreted as the description length of the game G, up to polynomial factors. In particular, for constant m, the game G can be described using  $|G|^{Op1q}$  bits. In Section 5, we discuss the interpretation of |G| when m is large.

The SparseMarkovCCE problem. From Fact C.1, we know that the problem of computing a sequence play:::;ptq of joint product Markov policies for which each player has at most T regret is equivalent to computing a sequence plq;:::;ptq for which the uniform mixture forms an -approximate CCE. We define pT;q-SparseMarkovCCE as the computational problem of computing such a CCE directly.

Definition 3.1 (SPARSEMARKOVCCE problem). For an m-player Markov game G and parameters T P N and i 0 (which may depend on the size of the game G), pT; q-SPARSEMARKOVCCE is the problem of finding a sequence  $^{\rm p1q}; \dots; ^{\rm pTq}$ , with each  $^{\rm ptg}$  P  $^{\rm markov}$ , such that the distributional policy  $^{\rm -1}$  T  $_{\rm t1}$  I  $_{\rm ptq}$  P  $^{\rm pgen;rnd}$ q is an -CCE of G (or equivalently, such that for all i P rms, Regi;T p^{\rm p1q}; \dots; ^{\rm pTq}q/T).

Decentralized learning algorithms naturally lead to solutions to the Sparse Markov CCE problem. In particular, consider any decentralized protocol which runs for T episodes, where at each timestep tPrTs, each player iPrms chooses a Markov policy ptq P markov to play, without knowledge of the other players' policies ptq (but possibly using the history); any strategy in which players independently run online learning algorithms falls under this protocol. If each player experiences overall regret at most T, then the sequence p1q;:::;pTq is a solution to the pT; q-SparseMarkovCCE problem. However, one might expect the pT;q-SparseMarkovCCE problem to be much easier than decentralized learning, since it allows for algorithms that produce p<sup>p1q</sup>;:::;<sup>pTq</sup>q satisfying the constraints of Definition 3.1 in a centralized manner. The main result of this section, Theorem 3.2, rules out the existence of any efficient algorithms, including centralized ones, that solve the Sparse Markov CCE problem.

Before moving on, let us give a sense for what sort of scaling one should expect for the parameters T and in the pT; q-SPARSEMARKOVCCE problem. First, we note that there always exists a solution to the p1;0q-SPARSEMARKOVCCE problem in a Markov game, which is given by a (Markov) Nash equilibrium of the game; of course, Nash equilibria are intractable to compute in general. For the special case of normal-form games (where there is only a single state, and H 1), no-regret learning (e.g., Hedge) yields a computationally efficient solution to the pT;Op1 Tqq-SPARSEMARKOVCCE problem, where the Opq hides a max<sub>i</sub>log |A<sub>i</sub>| factor. Refined convergence guarantees of Daskalakis et al. (2021); Anagnostides et al. (2022) improve upon this result, and yield an efficient solution

<sup>\*</sup>An alternative model allows for player i to have knowledge of the previous joint policies p1q;...;p1q, when selecting p1q.

 $<sup>^9</sup>$ In Appendix B, we discuss the implications of relaxing the stipulation that  $^{ptq}$  be product policies (for example, by allowing the use of shared randomness, as in V-learning). In short, allowing  $^{ptq}$  to be non-product essentially trivializes the problem.

<sup>&</sup>lt;sup>10</sup>We emphasize that pGq is defined as the maximum number of bits required by any particular ps;aq pair, not the total number of bits required for all ps;aq pairs.

<sup>&</sup>lt;sup>11</sup>Such a Nash equilibrium can be seen to exist by using backwards induction to specify the player's joint distribution of play at each state at steps H;H1;:::;1.

to the pT; Op1{Tqq-SparseMarkovCCE problem.

#### 3.2. Main result

Theorem 3.2. There is a constant  $C_0$  | 1 so that the following holds. Let nPN be given, and let TPN and | 0 satisfy T expp $^2$   $n^{1/2}$  { $2^5$ q. Suppose there is an algorithm that, given the description of any 2-player Markov game G with |G|/n, solves the pT;q-SparseMarkovCCE problem in time U, for some UPN. Then, for each nPN, the 2-player  $ptn^{1/2}u$ ;4q-Nash problem (Definition C.2) can be solved in time  $pnTUq^{C_0}$ .

We emphasize that the range T exppn<sup>Op1q</sup>q ruled out by Theorem 3.2 is the most natural parameter regime, since the runtime of any decentralized algorithm which runs for T episodes and produces a solution to the SparseMarkovCCE problem is at least linear in T. Using that 2-player pn;q-Nash is PPAD-complete for n<sup>c</sup> (for any c i 0) (Daskalakis et al., 2009; Chen et al., 2006; Rubinstein, 2018), we obtain the following corollary.

Corollary 3.3 (SPARSEMARKOVCCE is PPAD-complete). For any constant C  $_{\rm i}$  4, if there is an algorithm which, given the description of a 2-player Markov game G, solves the p|G| $^{\rm c}$ ; |G| $^{\rm c}$  $^{\rm d}$ -SPARSEMARKOVCCE problem in time polyp|G|q, then PPAD P.

The condition C i  $\frac{4}{|G|}$  in Corollary 3.3 is set to ensure that  $|G|^C$  expp $|G|^{2(C)}$   $\frac{4}{|G|}$  in Corollary 3.3 is set to ensure that to satisfy the condition of Theorem 3.2. Corollary 3.3 rules out the existence of a polynomial-time algorithm that solves the SparsemarkovCCE problem with accuracy polynomially small and T polynomially large in jGj.

Proof overview. The proof of Theorem 3.2 is based on a reduction, which shows that any algorithm that efficiently solves the pT;q-SparseMarkovCCE problem, for T not too large, can be used to efficiently compute an approximate Nash equilibrium of any given normal-form game. In particular, fix  $n_0 \, P \, N$ , and let a 2-player normal form game G with  $n_0$  actions be given.

We construct a Markov game G GpGq with horizon H n<sub>0</sub> and action sets identical to those of the game G, i.e., A<sub>1</sub> A<sub>2</sub> rn<sub>0</sub>s. The state space of G consists n<sup>2</sup> states, which are indexed by joint action profiles; the transitions are defined so that the value of the state at step h encodes the action profile taken by the agents at step h1.<sup>12</sup> At each state of G, the reward functions are given by the payoff matrices of G, scaled down by a factor of 1{H (which ensures that the rewards received at each step belong to r0;1{Hs}). In particular, the rewards and transitions out of a given state do not depend on the identity of the state, and so G can be thought of as a repeated game where G is played H times. The formal definition of G is given in Definition D.3.

Fix any algorithm for the SPARSEMARKOVCCE prob-

lem, and recall that for each step h and state s for G, ptq psq PpA1qpA2q denotes the joint action distribution taken in s at step h for the sequence of p1q;:::;pTq produced by the algorithm. The bulk of the proof of Theorem 3.2 consists of proving a key technical result, Lemma D.4, which states that if p1q;:::;pTq indeed solves pT;q-SparseMarkovCCE, then there exists some tuple ph;s;tq such that ptq psq is an approximate Nash equilibrium for G. With this established, it follows that we can find a Nash equilibrium efficiently by simply trying all HST choices for ph;s;tq.

To prove Lemma D.4, we reason as follows. Assume that  $\frac{1}{T} = \frac{1}{T} \frac{1}$ 

To sketch the idea, recall that to draw a trajectory from, we first draw a random index trTs uniformly at random, and then execute ptq for an episode. We will show (roughly) that for each player i, it is possible to compute a non-Markov deviation policy which, under the draw of a trajectory from can "infer" the value of the index t within the first few steps of the episode. Then policy then, at each state s and step h after the first few steps, play a best response to their opponent's portion of the strategy ptq psq. If, for each possible value of t, none of the distributions ptq psq are approximate Nash equilibria of G, this means that at least one of the players i can significantly increase their value in G over that of by playing; which contradicts the assumption that is an-CCE.

It remains to explain how we can construct a non-Markov policy  $\frac{1}{1}$ , which "infers" the value of t. Unfortunately, exactly inferring the value of t in the fashion described above is impossible: for instance, if there are  $t_1t_2$  so that  $pt_1q pt_2q$ , then clearly it is impossible to distinguish between the cases  $t t_1$  and  $t t_2$ . Nevertheless, by using the fact that each player ob-serves the full joint action profile played at each step h, we can construct a non-Markov policy which employs Vovk's aggregat-ing algorithm for online density estimation (Vovk, 1990; Cesa-

Bianchi & Lugosi, 2006) in order to compute a distribution which is close to people for most h PrHs. This guarantee is stated formally in an abstract setting in Proposition D.2, and is instantiated in the proof of Theorem 3.2 in ((5)). As we show in Section D.2, approximating people as we have described is sufficient to carry out the reasoning from the previous paragraph.

<sup>— &</sup>lt;sup>12</sup>For technical reasons, this only is the case for even values of h; we discuss further details in the full proof in Section D.2.

<sup>&</sup>lt;sup>13</sup>Vovk's aggregating algorithm is essentially the exponential weights algorithm with the logarithmic loss. A detailed background for the algorithm is provided in Section D.1.

#### 4. Lower bound for non-Markov algorithms

In this section, we prove Theorem 1.3 (restated formally below as Theorem 4.3), which strengthens Theorem 3.2 by allowing the sequence play;:::;prq of product policies to be non-Markovian. This additional strength comes at the cost of our lower bound only applying to 3-player Markov games (as opposed to Theorem 3.2, which applied to 2-player games).

#### 4.1. Sparsecce problem and computational model

To formalize the computational model for the SparseCCE problem, we must first describe how the non-Markov product policies ptq pptq;::;;ptq are represented. Recall that a non-Markov policy ptq Pgen;rnd is, by definition, a mapping from agent i's history and current state to a distribution over their next action. Since there are exponentially many possible histories, it is information-theoretically impossible to express an arbitrary policy in gen;rnd with polynomially many bits. As our focus is on computing a sequence of such policies ptq in polynomial time, certainly a prerequisite is that ptq can be expressed in polynomial space. Thus, we adopt the representational assumption, stated formally in Definition 4.1, that each of the policies  $_{i}^{ptq}$   $P_{i}^{gen;rnd}$  is described by a bounded-size circuit that can compute the conditional distribution of each next action given the history. This assumption is satisfied by essentially all empirical and theoretical work concerning non-Markov policies (e.g., (Leibo et al., 2021; Agapiou et al., 2022; Jin et al., 2021; Song et al., 2022)).

Definition 4.1 (Computable policy). Given a m-player Markov game G and N P N, we say that a policy  $_i$   $_i^{gen;rnd}$  is N-computable if for each h P r Hs, there is a circuit of size N that,  $_i^{14}$  on input  $_{p_i;h1}$ ; sq P  $_{h_i;h1}$  S, outputs the distribution  $_{ip_i;h1}$ ; sq P p A  $_{ip_i;h1}$ ; a policy  $_{p_i;h1}$ ;  $_{ip_i;h1}$ ;  $_{ip_i;h$ 

Our lower bound applies to algorithms that produce sequences  $p^{1q}; \dots; p^{Tq}$  for which each  $p^{tq}$  is N-computable, where the value N is taken to be polynomial in the description length of the game G. For example, Markov policies whose probabilities can be expressed with bits are OpHSA<sub>i</sub>q-computable for each player i, since one can simply store each of the probabil-ities  $p_{ij}$  ps<sub>h</sub>;  $p_{ij}$ , each of which takes bits to represent.

The Sparsecce problem. Sparsecce is the problem of computing a sequence of non-Markov product policies piq; :::; pTq such that the uniform mixture forms an approximate CCE. The problem generalizes Sparse-Markovcce (Definition 3.1) by relaxing the condition that the policies ptq be Markov.

Definition 4.2 (SPARSECCE Problem). For an m-player Markov game G and parameters T;N PN and i O (which may depend on the size of the game G), pT;;Nq-SPARSECCE is the problem of finding a sequence  $^{p1q}$ ;...; $^{pTq}$   $^{pgen;rnd}$ , with each  $^{ptq}$  being N-computable, such that the distributional policy  $^{-1}$ 

#### 4.2. Main result

Our main theorem for this section, Theorem 4.3, shows that for appropriate values of T, and N, solving the pT;;Nq-SPARSECCE problem is at least as hard as computing Nash equilibria in normal-form games.

Theorem 4.3. Fix n P N, and let T;N P N, and j 0 satisfy 1 T exp  $\frac{\hbar}{1-6}$ . Suppose there exists an algorithm that, given the description of any 3-player Markov game G with  $|G| \neq n$ , solves the pT;;Nq-SparseCCE problem in time U, for some U P N. Then, for any j 0, the 2-player ptn{2u; 50q-Nash problem can be solved in randomized time pnTNUlogp1{q{ $q^{C_0}$  with failure probability , where  $C_0$  j 0 is an absolute constant.

By analogy to Corollary 3.3, we obtain the following immediate consequence.

Corollary 4.4 (SPARSECCE is hard under PPAD  $\dagger$  RP). For any C ; 4, if there is an algorithm which, given the description of a 3-player Markov game G, solves the p|G|C;|G|LG|Cq-SPARSECCE problem in time polyp|G|Q, then PPAD, RP.

Proof overview for Theorem 4.3. The proof of Theorem 4.3 has a similar high-level structure to that of Theorem 3.2: given an m-player normal-form G, we define an pm 1q-player Markov game G GpGq which has no:tn{mu actions per player and horizon H no. The key difference in the proof of Theorem 4.3 is the structure of the players' reward functions. To motivate this difference and the addition of an pm 1q-th player, we explain why the proof of Theorem 3.2 fails to extend: a sequence p1q;::::p1q can hypothetically solve the Sparsecce problem by attempting to punish any one player's deviation policy, and thus avoid having to compute a Nash equilibrium of G. In particular, if player i plays according to the policy that we described in Section 3.2, then other players j i can use the non-Markov property of p1q to adjust their choice of actions in later rounds to decrease player i's value.

This behavior is reminiscent of "tit-for-tat" strategies which are used to establish the folk theorem in the theory of repeated games (Maskin & Fudenberg, 1986). The folk theorem describes how Nash equilibria are more numerous in repeated games than in single-shot normal form games. As it turns out, the folk theorem does not yield to worst-case speedups in repeated games, when the number of players is at least 3. Indeed, Borgs et al. (2008) gave an "anti-folk theorem", showing that computing Nash equilibria in pm 1q-player

<sup>&</sup>lt;sup>14</sup>For concreteness, we suppose that "circuit" means "boolean circuit" as in Definition 6.1 of (Arora & Barak, 2006), where probabilities are represented in binary. The precise model of computation we use does not matter, though, and we could equally assume that the policies i may be computed by Turing machines that terminate after N steps.

repeated games is PPAD-hard for m ¥ 2, via a reduction to m-player normal-form games. We adapt their reduction to our setting: roughly speaking, this approach adds an pm 1q-th player whose actions represent potential deviations for each of the m\_players. The structure of the rewards ensures that if  $\frac{1}{2} \int_{1}^{1} I_{ptq}$  is an -CCE, then for some policy  $\int_{1}^{1} I_{ptq}$  of the pm 1q-th player, the first m players will play an approximate Nash of G with constant probability, under a trajectory drawn from the joint policy  $\int_{1}^{1} I_{ptq} I_{p$ 

Two-player games. One intruiging question we leave open is whether the Sparsecce problem remains hard for two-player Markov games. Interestingly, as shown by Littman & Stone (2005), there is a polynomial time algorithm to find an exact Nash equilibrium for the special case of repeated two-player normal-form games. Though their result only applies in the infinite-horizon setting, it is possible to extend their results to the finite-horizon setting, which rules out naive approaches to extend the proof of Theorem 4.3 and Corollary 4.4 to two players.

#### 5. Multi-player games: lower bounds

In this section we present Theorem 1.4 (restated formally below as Theorem 5.2), which gives a statistical lower bound for the Sparsecce problem. The lower bound applies to any algorithm, regardless of computational cost, that accesses the underlying Markov game through a generative model.

Definition 5.1 (Generative model). For an m-player Markov game  $GpS;H;pA_iq_{iPrms};P;pR_iq_{iPrms};q$ , a generative model oracle is defined as follows: given a query described by a tuple ph;s;aqPrHsSA, the oracle returns the distribution  $P_hp|s;aqPpSq$  and the tuple of rewards  $pR_{i;h}ps;aqq_{iPrms}$ .

From the perspective of lower bounds, the assumption that the algorithm has access to a generative model is quite reasonable, as it encompasses most standard access models in RL, including the online access model, in which the algorithm repeatedly queries a policy and observes a trajectory drawn from it, as well as the local access generative model used in from (Yin et al., 2022; Weisz et al., 2021). We remark that it is slightly more standard to assume that queries to the generative model only return a sample from the distribution Php|s;aq as opposed to the distribution itself (Kakade, 2003; Kearns et al., 1999), but since our goal is to prove lower bounds, the notion in Definition 5.1 only makes our results stronger.

To state our main result, we recall the definition |G| maxtS; max<sub>iPrms</sub> A<sub>i</sub>; H; pGqu. In the present section, we consider the setting where the number of players m is large. Here, |G| does not necessarily correspond to the

description length for G, and should be interpreted, roughly speaking, as a measure of the description complexity of G [G] with respect to decentralized learning algorithms. In particular, from the perspective of an individual agent implementing a decentralized learning algorithm, their sample complexity should depend only on the size of their individual action set (as well as the global parameters S;H;pGq), as opposed to the size of the joint action set, which grows exponentially in m; the former is captured by jGj, while the latter is not. Indeed, a key advantage shared by much prior work on decentralized R L (Jin et al., 2021; Song et al., 2022; Mao & Basar, 2021; Daskalakis et al., 2022) is their avoidance of the curse of multi-agents, which describes the situation where an algorithm has sample and computational costs that scale exponentially in m.

Our main result for this section, Theorem 5.2, states that for m-player Markov games, exponentially many generative model queries (in m) are necessary to produce a solution to the pT;;Nq-SPARSECCE problem, unless T is exponential in m.

Theorem 5.2. Let m ¥ 2 be given. There are constants c; i 0 so that the following holds. Suppose there is an algorithm B which, given access to a generative model for a pm 1q-player Markov game G with  $|G| / 2m^6$ , solves the pT;{p10mq;Nq-SparseCCE problem for G for some T satisfying 1 T exppcmq, and any N P N. Then B must make at least 2  $p^{mq}$  queries to the generative model.

Theorem 5.2 establishes that there are m-player Markov games, where the number of states, actions per player, and horizon are bounded by polypmq, but any algorithm with regret opT {mq must make 2

queries (via Fact C.1). In particular, if there are polypmq queries per episode, as is standard in the online simulator model where a trajectory is drawn from the policy ptg at each episode t P r T s, then T j 2

stark contrast to the setting of normal-form games, where even for the case of bandit feedback (which is a special case of the generative model setting), stan-dard no-regret algorithms have the property that each player's regret scales as Op Tnq (i.e., independently of m), where n de-notes the number of actions per player (Lattimore & Szepesvari, 2020). As with our computational lower bounds, Theorem 5.2 is not limited to decentralized algorithms, and also rules out cen-tralized algorithms which have access to a generative model.

#### Acknowledgements

This work was performed in part while Noah Golowich was an intern at Microsoft Research. Noah Golowich is supported at MIT by a Fannie & John Hertz Foundation Fellowship and an NSF Graduate Fellowship. This work has been made possible in part by a gift from the Chan Zuckerberg Initiative Foundation to establish the Kempner Institute for the Study of Natural and Artificial Intelligence. Sham Kakade acknowledges funding from the Office of Naval Research under award N00014-22-1-2377 and the National Science Foundation Grant under award #CCF-2212841.

#### References

- Abbasi Yadkori, Y., Bartlett, P. L., Kanade, V., Seldin, Y., and Szepesvari, C. Online learning in markov decision processes with adversarially chosen transition probability distributions. In Burges, C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. (eds.), Advances in Neural Information Processing Systems, volume 26. Curran Associates, Inc., 2013. URL https://proceedings.neurips.cc/paper/2013/file/4f284803bd0966cc24fa8683a34afc6e-Paper.pdf.
- Agapiou, J. P., Vezhnevets, A. S., Duéñez-Guzmán, E. A., Matyas, J., Mao, Y., Sunehag, P., Köster, R., Madhushani, U., Kopparapu, K., Comanescu, R., Strouse, D., Johanson, M. B., Singh, S., Haas, J., Mordatch, I., Mobbs, D., and Leibo, J. Z. Melting pot 2.0, 2022. URL https://arxiv.org/abs/2211.13746.
- Anagnostides, I., Farina, G., Kroer, C., Lee, C.-W., Luo, H., and Sandholm, T. Uncoupled learning dynamics with \$o(nlog t)\$ swap regret in multiplayer games. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), Advances in Neural Information Processing Systems, 2022. URL <a href="https://openreview.net/forum?id=CZwh1XdAhNv">https://openreview.net/forum?id=CZwh1XdAhNv</a>.
- Arora, S. and Barak, B. Computational Complexity: A Modern Approach. Cambridge University Press, 2006. ISBN 978-0-521-42426-4.
- Babichenko, Y. Query complexity of approximate nash equilibria. J. ACM, 63(4), oct 2016. ISSN 0004-5411.
- Babichenko, Y. and Rubinstein, A. Communication complexity of approximate nash equilibria. In Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017, pp. 878–889, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450345286. doi: 10.1145/3055399.3055407. URL https://doi.org/10.1145/3055399.3055407.
- Bai, Y., Jin, C., and Yu, T. Near-optimal reinforcement learning with self-play. In Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Bakhtin, A., Brown, N., Dinan, E., Farina, G., Flaherty, C., Fried, D., Goff, A., Gray, J., Hu, H., Jacob, A. P., Komeili, M., Konath, K., Kwon, M., Lerer, A., Lewis, M., Miller, A. H., Mitts, S., Renduchintala, A., Roller, S., Rowe, D., Shi, W., Spisak, J., Wei, A., Wu, D., Zhang, H., and Zijlstra, M. Human-level play in the game of ¡i¿diplomacyj/i¿ by combining language models with strategic reasoning. Science, 378(6624):1067–1074, 2022. doi: 10.1126/science.ade9097. URL https://www.science.org/doi/abs/10.1126/science.ade9097.

- Blackwell, D. An analog of the minimax theorem for vector payoffs. Pacific Journal of Mathematics, 6:1–8, 1956.
- Borgs, C., Chayes, J., Immorlica, N., Kalai, A. T., Mirrokni, V., and Papadimitriou, C. The myth of the folk theorem. In Proceedings of the fortieth annual ACM symposium on Theory of computing, pp. 365–372, 2008.
- Brown, G. W. Some notes on computation of games solutions. 1949.
- Brown, N. and Sandholm, T. Superhuman ai for heads-up no-limit poker: Libratus beats top professionals. Science, 359(6374):418–424, 2018. doi: 10.1126/science.aao1733. URL https://www.science.org/doi/abs/10.1126/science.aao1733.
- Cesa-Bianchi, N. and Lugosi, G. Prediction, Learning, and Games. Cambridge University Press, New York, NY, USA, 2006. ISBN 0521841089.
- Cesa-Bianchi, N., Freund, Y., Haussler, D., Helmbold, D. P., Schapire, R. E., and Warmuth, M. K. How to use expert advice. Journal of the ACM, 44(3):427–485, 1997.
- Chen, X. and Peng, B. Hedging in games: Faster convergence of external and swap regrets. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 18990–18999. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/file/db346ccb62d491029b590bbbf0f5c412-Paper.pdf.
- Chen, X., Deng, X., and Teng, S.-h. Computing nash equilibria: Approximation and smoothed complexity. In 2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06), pp. 603–612, 2006.
- Chen, X., Cheng, Y., and Tang, B. Well-Supported vs. Approximate Nash Equilibria: Query Complexity of Large Games. In Papadimitriou, C. H. (ed.), 8th Innovations in Theoretical Computer Science Conference (ITCS 2017), volume 67 of Leibniz International Proceedings in Informatics (LIPIcs), pp. 57:1–57:9, Dagstuhl, Germany, 2017. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-029-3. doi: 10.4230/LIPIcs.ITCS.2017.57. URL http://drops.dagstuhl.de/opus/volltexte/2017/8163.
- Daskalakis, C., Goldberg, P. W., and Papadimitriou, C. H. The complexity of computing a nash equilibrium. SIAM Journal on Computing, 39(1):195–259, 2009.
- Daskalakis, C., Golowich, N., and Zhang, K. The complexity of markov equilibrium in stochastic games, 2022. URL https://arxiv.org/abs/2204.03991.
- Daskalakis, C. C., Fishelson, M., and Golowich, N. Near-optimal no-regret learning in general

- games. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), Advances in Neu-ral Information Processing Systems, 2021. URL https://openreview.net/forum?id=cVwc7IHWEWi.
- Erez, L., Lancewicki, T., Sherman, U., Koren, T., and Mansour, Y. Regret minimization and convergence to equilibria in general-sum markov games, 2022. URL https://arxiv.org/abs/2207.14211.
- Fearnley, J., Gairing, M., Goldberg, P., and Savani, R. Learning equilibria of games via payoff queries. In Proceedings of the Fourteenth ACM Conference on Electronic Commerce, EC '13, pp. 397–414, New York, NY, USA, 2013. Association for Computing Machinery. ISBN 9781450319621. doi: 10.1145/2492002.2482558. URL https://doi.org/10.1145/2492002.2482558.
- Foster, D. J., Rakhlin, A., Sekhari, A., and Sridharan, K. On the complexity of adversarial decision making. arXiv preprint arXiv:2206.13063, 2022.
- Fudenberg, D., Levine, D., and Maskin, E. The folk theorem with imperfect public information. Econometrica, 62(5): 997–1039, 1994. ISSN 00129682, 14680262.
- Hannan, J. Approximation to Bayes risk in repeated play. Contributions to the Theory of Games, 3:97–139, 1957.
- Hart, S. and Mas-Colell, A. A simple adaptive procedure leading to correlated equilibrium. Econometrica, 68(5):1127–1150, 2000. doi: https://doi.org/10.1111/1468-0262.00153. URL https://onlinelibrary.wiley.com/doi/abs/10.1111/1468-0262.00153.
- Jin, C., Krishnamurthy, A., Simchowitz, M., and Yu, T. Reward-free exploration for reinforcement learning. In International Conference on Machine Learning (ICML), pp. 4870–4879. PMLR, 2020.
- Jin, C., Liu, Q., Wang, Y., and Yu, T. V-learning—a simple, efficient, decentralized algorithm for multiagent RL. arXiv preprint arXiv:2110.14555, 2021.
- Jin, Y., Muthukumar, V., and Sidford, A. The complexity of infinite-horizon general-sum stochastic games, 2022.
- Kakade, S. M. On the sample complexity of reinforcement learning, 2003.
- Kearns, M., Mansour, Y., and Ng, A. Y. A sparse sampling algorithm for near-optimal planning in large markov decision processes. In Proceedings of the 16th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI'99, pp. 1324–1331, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.

- Kramár, J., Eccles, T., Gemp, I., Tacchetti, A., McKee, K. R., Malinowski, M., Graepel, T., and Bachrach, Y. Negotiation and honesty in artificial intelligence methods for the board game of Diplomacy. Nature Communications, 13(1):7214, December 2022. ISSN 2041-1723. doi: 10.1038/s41467-022-34473-5. URL https://www.nature.com/articles/s41467-022-
  - 34473-5. Number: 1 Publisher: Nature Publishing Group.
- Kwon, J., Efroni, Y., Caramanis, C., and Mannor, S. RL for latent MDPs: Regret guarantees and a lower bound. Advances in Neural Information Processing Systems, 34, 2021.
- Lattimore, T. and Szepesvári, C. Bandit algorithms. Cambridge University Press, 2020.
- Leibo, J. Z., Dueñez-Guzman, E. A., Vezhnevets, A., Agapiou, J. P., Sunehag, P., Koster, R., Matyas, J., Beattie, C., Mordatch, I., and Graepel, T. Scalable evaluation of multi-agent reinforcement learning with melting pot. In Meila, M. and Zhang, T. (eds.), Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 6187–6199. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/leibo21a.html.
- Littlestone, N. and Warmuth, M. K. The weighted majority algorithm. Information and computation, 108(2):212–261, 1994.
- Littman, M. L. and Stone, P. A polynomial-time Nash equilibrium algorithm for repeated games. Decision Support Systems, 39:55–66, 2005.
- Liu, Q., Wang, Y., and Jin, C. Learning Markov games with adversarial opponents: Efficient algorithms and fundamental limits. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), Proceedings of the 39th International Conference on Machine Learning, volume 162 of Proceedings of Machine Learning Research, pp. 14036–14053. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/liu22r.html.
- Malialis, K. and Kudenko, D. Distributed response to network intrusions using multiagent reinforcement learning. Engineering Applications of Artificial Intelligence, 41:270–284, 2015. ISSN 0952-1976. doi: https://doi.org/10.1016/j.engappai.2015.01.013. URL https://www.sciencedirect.com/science/article/pii/S095219761500024X.
- Mao, W. and Basar, T. Provably efficient reinforcement learning in decentralized general-sum markov games. CoRR, abs/2110.05682, 2021. URL https://arxiv.org/abs/2110.05682.
- Maskin, E. and Fudenberg, D. The folk theorem in repeated games with discounting or with incomplete information.

- Econometrica, 53(3):533–554, 1986. Reprinted in A. Rubinstein (ed.), Game Theory in Economics, London: Edward Elgar, 1995. Also reprinted in D. Fudenberg and D. Levine (eds.), A Long-Run Collaboration on Games with Long-Run Patient Players, World Scientific Publishers, 2009, pp. 209-230.
- Nash, J. Non-cooperative games. Annals of Mathematics, 54 (2):286–295, 1951. ISSN 0003486X.
- Papadimitriou, C. H. On the complexity of the parity argument and other inefficient proofs of existence. J. Comput. Syst. Sci., 48(3):498–532, 1994.
- Perolat, J., Vylder, B. D., Hennes, D., Tarassov, E., Strub, F., de Boer, V., Muller, P., Connor, J. T., Burch, N., Anthony, T., McAleer, S., Elie, R., Cen, S. H., Wang, Z., Gruslys, A., Malysheva, A., Khan, M., Ozair, S., Timbers, F., Pohlen, T., Eccles, T., Rowland, M., Lanctot, M., Lespiau, J.-B., Piot, B., Omidshafiei, S., Lockhart, E., Sifre, L., Beauguerlange, N., Munos, R., Silver, D., Singh, S., Hassabis, D., and Tuyls, K. Mastering the game of stratego with model-free multiagent reinforcement learning. Science, 378(6623):990–996, 2022. doi: 10.1126/science.add4679. URL https://www.science.org/doi/abs/10.1126/science.add4679.
- Puterman, M. Markov Decision Processes. John Wiley & Sons, Ltd, 1 edition, 1994. URL http://onlinelibrary.wiley.com/doi/10.1002/9780470316887.
- Roughgarden, T. Intrinsic robustness of the price of anarchy. Journal of the ACM, 2015.
- Rubinstein, A. Settling the complexity of computing approximate two-player Nash equilibria. In Annual Symposium on Foundations of Computer Science (FOCS), pp. 258–265. IEEE, 2016.
- Rubinstein, A. Inapproximability of Nash equilibrium. SIAM Journal on Computing, 47(3):917–959, 2018.
- Shalev-Shwartz, S., Shammah, S., and Shashua, A. Safe, multi-agent, reinforcement learning for autonomous driving. CoRR, abs/1610.03295, 2016. URL http://arxiv.org/abs/1610.03295.
- Shapley, L. Stochastic Games. PNAS, 1953.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. nature, 529(7587):484, 2016.
- Song, Z., Mei, S., and Bai, Y. When can we learn general-sum markov games with a large number of players sample-efficiently? In International Conference on Learning Representations, 2022. URL https://openreview.net/forum?id=6MmiS0HUJHR.

- Syrgkanis, V., Agarwal, A., Luo, H., and Schapire, R. E. Fast convergence of regularized learning in games. In Advances in Neural Information Processing Systems (NIPS), pp. 2989–2997, 2015.
- Vovk, V. Aggregating strategies. Proc. of Computational Learning Theory, 1990, 1990.
- Weisz, G., Amortila, P., Janzer, B., Abbasi-Yadkori, Y., Jiang, N., and Szepesvari, C. On query-efficient planning in mdps under linear realizability of the optimal state-value function. In Belkin, M. and Kpotufe, S. (eds.), Proceedings of Thirty Fourth Conference on Learning Theory, volume 134 of Proceedings of Machine Learning Research, pp. 4355–4385. PMLR, 15–19 Aug 2021.
- Yin, D., Hao, B., Abbasi-Yadkori, Y., Lazić, N., and Szepesvári, C. Efficient local planning with linear function approximation. In Dasgupta, S. and Haghtalab, N. (eds.), Proceedings of The 33rd International Conference on Algorithmic Learning Theory, volume 167 of Proceedings of Machine Learning Research, pp. 1165–1192. PMLR, 29 Mar–01 Apr 2022.
- Zhan, W., Lee, J. D., and Yang, Z. Decentralized optimistic hyperpolicy mirror descent: Provably no-regret learning in markov games. arXiv preprint arXiv:2206.01588, 2022.
- Zheng, S., Trott, A., Srinivasa, S., Parkes, D. C., and Socher, R. The ai economist: Taxation policy design via two-level deep multiagent reinforcement learning. Science Advances, 8(18):eabk2607, 2022. doi: 10.1126/sciadv.abk2607. URL https://www.science.org/doi/abs/10.1126/sciadv.abk2607.

### **Contents of Appendix**

I	Additional results and discussion	14
Α	Tighter computational lower bounds under ETH for <b>PPAD</b>	14
В	Discussion and interpretation	15
	B.1 Comparison to V-learning	B.2 No-regre
	learning against Markov deviations	e role of shared
	randomness	wer bounds fo
	finding stationary CCE	
П	Proofs	17
С	Additional preliminaries	18
	C.1 Additional preliminaries for Markov games	18 C.2 Nash
	equilibria and computational hardness	ery complexity
	of Nash equilibria	
D	Proofs of lower bounds for Sparse MarkovCCE (Section 3)	19
	D.1 Preliminaries: Online density estimation	9 D.2 Proof o
	Theorem 3.2	
E	Proofs of lower bounds for SparseCCE (Sections 4 and 5)	24
	E.1 Proof of Theorem E.1	E.2 Remarks
	on bit complexity of the rewards	
1	F Equivalence between $_{j}^{\text{gen;rnd}}$ and $_{j}^{\text{gen;det}}$ q	33 F.1
	Proofs of the equivalence	34

#### Part I

## Additional results and discussion

#### A. Tighter computational lower bounds under ETH for PPAD

Recall that Corollary 3.3 states that if PPAD P, then there is no constant C  $_{\rm i}$  4 and polyp $_{\rm i}$ q-time algorithm which solves the p $_{\rm i}$ q-SparseMarkovCCE problem for any 2-player Markov game G. Using a stronger complexity-theoretic assumption, the Exponential Time Hypothesis for PPAD (Rubinstein, 2016), we can obtain a hardness result which rules out

efficient algorithms even when 1) the accuracy is constant, as opposed to being  $|G|^{1/C}$ , and 2) T is quasipolynomially large, as opposed to only being of polynomial size, i.e.,  $|G|^{C}$ .

Corollary A.1 (ETH-hardness of SparseMarkovCCE). There is a constant  $_0$  i  $_0$  such that if there exists an algorithm that solves the  $p|G|^{oplog|G|q}$ ;  $_0q$ -SparseMarkovCCE problem in  $|G|^{oplog|G|q}$  time, then the Exponential Time Hypothesis for PPAD fails to hold.

Corollary A.1 is an immediate consequence of Theorem 3.2 and the fact that for some absolute constant  $_0$  i 0, there are no polynomial-time algorithms for computing  $_0$ -Nash equilibria in 2-player normal-form games under the Exponential Time Hypothesis for PPAD (as shown in (Rubinstein, 2016)).

#### B. Discussion and interpretation

Theorems 3.2, 4.3, and 5.2 present barriers—both computational and statistical—toward developing efficient decentralized no-regret guarantees for multi-agent reinforcement learning. We emphasize that no-regret algorithms are the only known approach for obtaining fully decentralized learning algorithms (i.e., those which do not rely even on shared randomness) in normal-form games, and it seems unlikely that a substantially different approach would work in Markov games. Thus, these lower bounds for finding subexponential-length sequences of policies with the no-regret property represent a significant obstacle for fully decentralized multi-agent reinforcement learning. Moreover, these results rule out even the prospect of developing efficient centralized algorithms that produce no-regret sequences of policies, i.e., those which "resemble" independent learning. In this section, we compare our lower bounds with recent upper bounds for decentralized learning in Markov games, and explain how to reconcile these results.

#### B.1. Comparison to V-learning

The V-learning algorithm (Jin et al., 2021; Song et al., 2022; Mao & Basar, 2021) is a polynomial-time decentralized learning algorithm that proceeds in two phases. In the first phase, the m agents interact over the course of K episodes in a decentralized fashion, playing product Markov policies parking P markov. In the second phase, the agents use data gathered during the first phase to produce a distributional policy p P pgen;rndq, which we refer to as the output policy of V-learning. As discussed in Section 1, one implication of Theorem 3.2 is that the first phase of V-learning cannot guarantee each agent sublinear regret. Indeed if K is of polynomial size (and PPAD P), this follows because a bound of the form Reg i; K ppiq;:::;pkqq / K for all i implies that ppiq;:::;pkqq solves the pK;q-SparseMarkovCCE problem.

The output policy  $\beta p^{gen;rnd}q$  produced by V-learning is an approximate CCE (per Definition 2.1), and it is natural to ask how many product policies it takes to represent as a uniform mixture (that is, whether solves the pT;q-SPARSEMARKOVCCE problem for a reasonable value of T). First, recall that V-learning requires K polypH;S;max<sub>i</sub>A<sub>i</sub>q $\{^2$  episodes to ensure that p is an -CCE. It is straightforward to show that p can be expressed as a non-uniform mixture of at most K <sup>K H S 1</sup> policies in <sup>gen;rnd</sup> (we prove this fact in detail below). By discretizing the non-uniform mixture, one can equivalently represent it as uniform mixture of Op1{q K <sup>KHS 1</sup> product policies, up to error. Recalling the value of K, we conclude that we can express as p a uniform mixture of T exppOp1{2qpolypH;S;max<sub>i</sub>A<sub>i</sub>qq product policies in <sup>gen;rnd</sup>. Note that the lower bound of Theorem 4.3 rules out the efficient computation of an -CCE represented as a uniform mixture of T! expp<sup>2</sup> maxtH;S; max<sub>i</sub>A<sub>i</sub>uq efficiently computable policies in <sup>gen;rnd</sup>. Thus, in the regime where 1{ is polynomial in H;S; max<sub>i</sub>A<sub>i</sub>, this upper bound on the sparsity of the policy p produced V-learning matches that from Theorem 4.3, up to a polynomial in the exponent.

The sparsity of the output policy from **V-learning**. We now sketch a proof of the fact that the output policy produced by V-learning can be expressed as a (non-uniform) average of K <sup>K H S 1</sup> policies in <sup>gen;rnd</sup>, where K is the number of episodes in the algorithm's initial phase. We adopt the notation and terminology from Jin et al. (2021).

Consider Algorithm 3 of Jin et al. (2021), which describes the second phase of V-learning, which produces the output policy p. We describe how to write p as a weighted average of a collection of product policies, each of which is indexed by a function : rHsSrKsÑrKs and a parameter  $k_0$ PrKs: in particular, we will write p  $k_0$ ;  $w_{k_0}$ 

We define the mixing weight allocated  $w_{k_0}$ ; to any tuple  $pk_0$ ;q to be:

$$\frac{1}{K} \underset{ph;s;kqPrHsSrKs}{"} 1 tph;s;kqPrN^k psqsu^{ph;s;kq} \underset{N}{\overset{}{}_{\not k}} p_{psq}$$

where N  $_h^k$ psq PrKs and  $_h^{Nips}$ q PrO;1s (for iPrN $_h$ p/sqs) are defined as in (Jin et al., 2021).

Next, for each  $k_0$ ;, we define  $_k$ ;  $P_0^{gen;rnd}$  to be the following policy: it maintains a parameter k P r Ks over the first h H steps of the episode (as in Algorithm 3 of (Jin et al., 2021)), but upon reaching state s at step h, given the present value of k P r Ks, sets i:ph;s;kq, and updates k —  $k^i$  psq, and then samples an action  $a^kp$  sq (where  $k^i$  psq,  $k^i$  psq defined in (Jin et al., 2021)). Since the mixing weights  $w_k$ ; defined above exactly simulate the random draws of the parameter k in Line 1 and the parameters k in Algorithm 3, Line 4 of (Jin et al., 2021), is equal to  $k_0; w_{k_0}; k_0; P p^{gen;rnd}q.$ 

#### B.2. No-regret learning against Markov deviations

As discussed in Section 1, Erez et al. (2022) showed the existence of a learning algorithm with the property that if each agent plays it independently for T episodes, then no player can achieve regret more than Oppolypm;H;S;max<sub>i</sub>A<sub>i</sub>qT<sup>3{4</sup>q by deviating to any fixed Markov policy. This notion of regret corresponds to, in the context of Definition 2.2, replacing max penint with the smaller quantity max pmarkov. Thus, the result of Erez et al. (2022) applies to a weaker notion of regret than that of the Sparsece problem, and so does not contradict any of our lower bounds. One may wonder which of these two notions of regret (namely, best possible gain via deviation to a Markov versus non-Markov policy) is the "right" one. We do not believe that there is a definitive answer to this question, but we remark that in many empirical applications of multi-agent reinforcement learning it is standard to consider non-Markov policies (Leibo et al., 2021; Agapiou et al., 2022). Furthermore, as shown in the proposition below, there are extremely simple games, e.g., of constant size, in which Markov deviations lead to "vacuous" behavior: in particular, all Markov policies have the same (suboptimal) value but the best non-Markov policy has much greater value:

Proposition B.1. There is a 2-player, 2-action, 1-state Markov game with horizon 2 and a non-Markov policy  $P_2^{\text{gen;rnd}}$  for player 2 so that for all  $_1P_1^{\text{markov}}$ ,  $V^{12}$   $_1P_2^{\text{markov}}$ ,  $V^{12}$   $_1P_2^{\text{gen;rnd}}$   $V^{12}$   $_1P_2^{\text{gen;rnd}}$   $V^{12}$   $_1P_2^{\text{markov}}$ ,  $V^{12}$ 

The proof of Proposition B.1 is provided in Section B.5 below.

Other recent work has also proved no-regret guarantees with respect to deviations to restricted policy classes. In particular, Zhan et al. (2022) studies a setting in which each agent i is allowed to play policies in an arbitrary restricted policy class  $_{i}^{1}$   $_{i}^{gen;rnd}$  in each episode, and regret is measured with respect to deviations to any policy in  $_{i}^{1}$ . Zhan et al. (2022) introduces an algorithm, DORIS, with the property that when all agents play it independently, each agent i experiences regret O polypm;A;S;Hq  $_{i}^{a}$   $_{i}^{a}$   $_{i}^{a}$   $_{i}^{a}$   $_{i}^{b}$  of their respective class  $_{i}^{1}$ .

DORIS is not computationally efficient, since it involves performing exponential weights over the class  $^1$ , which requires space complexity  $j^1j$ . Nonetheless, one can compare the statistical guarantees the algorithm provides to our own results. Let  $^{\text{markov}; \text{det}}$   $\mathfrak{E}^{\text{markov}}_{i}$  denote the set of deterministic Markov policies of agent i, namely sequences  $_i$   $p_{i;1}; \dots;_{i;H}q$  so that  $_{i;h}: S \ \tilde{N} \ A_i$ . In the case that  $_{\text{markov}; \text{det}}^1$ ,  $_{\text$ 

Markov deviations when m is constant, comparable to Erez et al. (2022). However, we are interested in the setting in which each player's regret is measured with respect to all deviations in  $_{i}^{\text{gen;rnd}}$  (equivalently,  $_{i}^{\text{gen;det}}$ ). Accordingly, if we take  $_{i}^{\text{gen;det}}$  equivalently,  $_{i}^{\text{gen;det}}$ ). Accordingly, if we take  $_{i}^{\text{gen;det}}$  then  $\log |_{i}^{\text{gen;det}}|_{i}^{\text{gen;det}}$ , meaning that DORIS does not imply any sort of sample-efficient guarantee, even for m2.

Finally, we remark that the algorithm DORIS (Zhan et al., 2022), as well as the similar algorithm OPMD from earlier work of Liu et al. (2022), obtains the same regret bound stated above even when the opponents are controlled by (possibly adaptive) adversaries. However, this guarantee crucially relies on the fact that any agent implementing DORIS must observe the policies played by opponents following each episode; this feature is the reason that the regret bound of DORIS does not contradict the exponential lower bound of Liu et al. (2022) for no-regret learning against an adversarial opponent. As a result of being restricted to this "revealed-policy" setting, DORIS is not a fully decentralized algorithm in the sense we consider in this paper.

<sup>&</sup>lt;sup>16</sup>Erez et al. (2022) has the added bonus of computational efficiency, even for polynomially large m, though has the significant drawback of assuming that the Markov game is known.

 $<sup>^{17} \</sup>text{DORIS}$  plays distributions over policies in  $^1$   $_{i}^{\text{gen;det}}$  at each episode, whereas in our lower bounds we consider the setting where a policy in  $^{\text{gen;rnd}}$  is played each episode; Facts F.2 and F.3 shows that these two settings are essentially equivalent, in that any policy in  $^{\text{gen;rnd}}$  can be simulated by one in  $p^{\text{gen;det}}qp^{\text{gen;det}}q$ , and vise versa.  $_{\text{m}}$ 

#### B.3. On the role of shared randomness

A key assumption in our lower bounds for no-regret learning is that each of the joint policies pta;::::,pta produced by the algorithm is a product policy; such an assumption is natural, since it subsumes independent learning protocols in which each agent i selects pta without knowledge of pta. Compared to general (stochastic) joint policies, product policies have the desirable property that, to sample a trajectory from pta ppta,::::;pta q Pen;rnd gen;rnd gen

#### B.4. Comparison to lower bounds for finding stationary CCE

A separate line of work Daskalakis et al. (2022); Jin et al. (2022) has recently shown PPAD-hardness for the problem of finding stationary Markov CCE in infinite-horizon discounted stochastic games. These results are incomparable with our own: stationary Markov CCE are not sparse (in the sense of Definition 3.1), whereas we do not require stationarity of policies (as is standard in the finite-horizon setting).

#### B.5. Proof of Proposition B.1

Below we prove Proposition B.1.

Proof of Proposition B.1. We construct the claimed Markov game G as follows. The single state is denoted by s; as there is only a single state, the transitions are trivial. We denote each player's action space as  $A_1 A_2 t1$ ;2u. The rewards to player 1 are given as follows: for all  $pa_1$ ; $a_2qPA$ ,

$$\mathsf{R}_{1;1}\mathsf{ps};\mathsf{pa}_1;\mathsf{a}_2\mathsf{qq} \ _2^{\ \ l}_{\mathsf{a}_2 1}; \qquad \qquad \mathsf{R}_{1;2}\mathsf{ps};\mathsf{pa}_1;\mathsf{a}_2\mathsf{qq} \ _2^{\ \ l}_{\mathsf{a}_1 \mathsf{a}_2}^{\ \ 1} :$$

We allow the rewards of player 2 to be arbitrary; they do not affect the proof in any way.

We let  ${}_2p_{2;1};_{2;2}qP^{gen;rnd}$  by the policy which plays a uniformly random action at step 1 and then plays the same action at step 2: formally,  ${}_{2;1}ps_1qUnifpA_2q$ , and  ${}_{2;2}pps_1;_{a_2;1};_{r_2;1}q;_{s_2}qI_{a_{2;1}}$ . Then for any Markov policy  ${}_1P^{markov}$  of player 1, we must have  $P_{12}pa_{1;2}$   $a_{2;2}q1$ {2, which means that  $V_1$   ${}^1$   ${}^2$   ${}^1$   ${}^2$   ${}^1$   ${}^2$   ${}^1$   ${}^2$   ${}^1$   ${}^2$   ${}^2$   ${}^2$   ${}^1$   ${}^2$ 

On the other hand, any general (non-Markov) policy  $_1$  P  $_1^{\rm gen;rnd}$  which satisfies

has 
$$V_1^{12}$$
 1{2p1{2 1q3{4.}}

#### Part II

## **Proofs**

#### C. Additional preliminaries

#### C.1. Additional preliminaries for Markov games

Additional facts on regret and CCE. The following facts regarding deterministic policies and the definition of coarse correlated equilibria and regret are well-known:

- In the context of (1) in the definition of regret, the maximum over  $_i$   $P_i^{gen;rnd}$  is always achieved by a deterministic general  $_{ptq}$  policy, so we have  $Reg_{i;T}$   $max_i$   $_{Pgen;det}$   $_{t1}$   $V_{iT}$   $_{i}$   $_{i}$   $_{i}$   $_{i}$   $V_{i}$   $_{ptq}$  .

Next, the following standard result shows that the uniform average of any no-regret sequence forms an approximate coarse correlated equilibrium.

Fact C.1 (No-regret is equivalent to CCE). Suppose that a sequence of policies  $^{p1q}$ ; :::;  $^{pTq}$  P  $^{gen;rnd}$  satisfies  $Reg_{i;T}$   $p^{p1q}$ ; :::;  $^{pTq}$ q  $^{\prime}$ . T for each i P rms. Then the uniform average of these T policies, namely the distributional policy:  $^{T}$   $^{T}$ 

Likewise if a sequence of policies  $^{p1q}$ ; ::: ;  $^{pTq}$  P  $^{gen;rnd}$  has the property that the distributional policy:  $^{T}$   $^{1}$   $^{T}$   $^{1}$   $^{$ 

Fact C.1 is an immediate consequence of Definitions 2.1 and 2.2.

#### C.2. Nash equilibria and computational hardness.

The most foundational and well known solution concept for normal-form games is the Nash equilibrium (Nash, 1951).

Definition C.2 (pn; q-Nash problem). For a normal-form game G  $pM_1$ ; :::;  $M_mq$  and j 0, a product distribution pP  $\frac{m}{i1}$ prnsq is said to be an -Nash equilibrium for G if for all iPrns,

$$\max_{a_i Prns} E_{ap} rpM_i q_{a_i;a_i} sE_{ap} rpM_i q_s /:$$

We define the m-player pn;q-NASH problem to be the problem of computing an -Nash equilibrium of a given m-player n-action normal-form game. 18

Informally, p is an -Nash equilibrium if no player i can gain more than in reward by deviating to a single fixed action  $a^1$ , while all other players randomly choose their actions according to p. Despite the intuitive appeal of Nash equilibria, they are intractable to compute: for any c i 0, it is PPAD-hard to solve the pn; $n^cq$ -NASH problem, namely, to compute  $n^c$ -approximate Nash

equilibria in 2-player n-action normal-form games (Daskalakis et al., 2009; Chen et al., 2006; Rubinstein, 2018). We recall that the complexity class PPAD consists of all total search problems which have a polynomial-time reduction to the End-of-The-Line (EOTL) problem. PPAD is the most well-studied complexity class in algorithmic game theory, and it is widely believed that PPAD P. We refer the reader to (Daskalakis et al., 2009; Chen et al., 2006; Rubinstein, 2018; Papadimitriou, 1994) for further background on the class PPAD and the EOTL problem.

#### C.3. Query complexity of Nash equilibria

Our statistical lower bound for the SPARSECCE problem in Theorem 5.2 relies on existing query complexity lower bounds for computing approximate Nash equilibria in m-player normal-form games. We first review the query complexity model for normal-form games.

Oracle model for normal-form games. For m;nPN, consider an m-player n-action normal form game G, specified by payoff tensors  $M_1;...;M_m$ . Since the tensors  $M_1;...;M_m$  contain a total of mn<sup>m</sup> real-valued payoffs, in the setting when m is large, it is unrealistic to assume that an algorithm is given the full payoff tensors as input. Therefore, prior work on computing equilibria in such games has studied the setting in which the algorithm makes adaptive oracle queries to the payoff tensors.

In particular, the algorithm, which is allowed to be randomized, has access to a payoff oracle  $O_G$  for the game G, which works as follows. At each time step, the algorithm can choose to specify an action profile a P rns<sup>m</sup> and then query  $O_G$  at the action profile a. The oracle  $O_G$  then returns the payoffs  $pM_1q_a$ ;:::; $pM_mq_a$  for each player if the action profile a is played.

Query complexity lower bound for approximate Nash equilibrium. The following theorem gives a lower bound on the number of queries any randomized algorithm needs to make to compute an approximate Nash equilibrium in an m-player game.

Theorem C.3 (Corollary 4.5 of (Rubinstein, 2016)). There is a constant  $_{0\,i}$  0 so that any randomized algorithm which solves the p2; $_{0}$ q-Nash problem for m-player normal-form games with probability at least 2{3 must use at least 2}  $_{pmq}$  payoff queries.

We remark that (Babichenko, 2016; Chen et al., 2017) provide similar, though quantitatively weaker, lower bounds to that in Theorem C.3. We also emphasize that the lower bound of Theorem C.3 applies to any algorithm, i.e., including those which require extremely large computation time.

#### D. Proofs of lower bounds for Sp<sup>A</sup> RSEMARKOVCCE (Section 3)

#### D.1. Preliminaries: Online density estimation

Our proof makes use of tools for online learning with the logarithmic loss, also known as conditional density estimation. In particular, we use a variant of the exponential weights algorithm known as Vovk's aggregating algorithm in the context of density estimation (Vovk, 1990; Cesa-Bianchi & Lugosi, 2006). We consider the following setting with two players, a Learner and Nature. Furthermore, there is a set Y, called the outcome space, and a set X, called the context space; for our applications it suffices to assume Y and X are finite. For some T PN, there are T time steps t1;2;:::;T. At each time step tPrTs:

- Nature reveals a context x<sup>ptq</sup> P X;
- Having seen the context  $x^{ptq}$ , the learner predicts a distribution  $\mathbf{p}^{pq}$  PpYq;
- Nature chooses an outcome  $y^{ptq}$  P Y , and the learner suffers loss  $'_{log}^{ptq} p p^{rtq} q : log q p q q = 0$ :

For each tPrTs, we let  $H^{ptq}$  tpx $^{p1q}$ ; $y^{p1q}$ ; $q^{p1q}$ ; $y^{ptq}$ ; $q^{ptq}$ ;

Note that the learner can observe the expert predictions  $tp_ipx^{ptq}qu_{iPl}$  and use them to make its own prediction at each round t.

Proposition D.1 (Vovk's aggregating algorithm). Consider Vovk's aggregating algorithm, which predicts via

This algorithm guarantees a regret bound of  $Reg_{I \cap T} / log |I|$ .

Recall that for probability distributions p;q on a finite set B, their total variation distance is defined as

$$D_{TV}pp;qq \max_{E \notin R} |ppEqqpEq|$$
: (3)

As a (standard) consequence of Proposition D.1, in the realizable setting in which the distribution of y<sup>ptq</sup> | x<sup>ptq</sup> follows p<sub>i</sub>px<sup>ptq</sup> for some fixed (unknown) expert i PI, we can obtain a bound on the total variation distance between the algorithm's predictions and those of  $p_i p x^{ptq} q$ .

Proposition D.2. If the distribution of outcomes is realizable, i.e., there exists an expert i P I so that  $y^{ptq}$   $p_i p x^{ptq}$   $p_i p x^{ptq}$ t PrTs, then the predictions p tq of the aggregation algorithm (2) satisfy

$$\begin{array}{c}
T \\
\text{ErD}_{\text{TV}} p \mathbf{p} p^{\text{ptq}}; p_i p x^{\text{ptq}} q q s / \underline{a} \\
T \log | | : \\
t1
\end{array}$$

For completeness, we provide the proof of Proposition D.2 here.

Proof of Proposition D.2. To simplify notation, for an expert i P I, a context x P X, and an outcome y P Y, we write p<sub>i</sub>py|xq to denote pipxqpyq.

Proposition D.1 gives that the following inequality holds (almost surely):

For each tPrTs, note that  $\mathbf{p}^{\text{tq}}$  and  $\mathbf{x}^{\text{ptq}}$  are  $\mathbf{F}^{\text{pt1q}}$ -measurable (by definition). Then

where the first inequality uses Pinsker's inequality and the final equality uses the fact that  $y^{ptq} p_i p x^{ptq} q | x^{ptq} ; H^{pttq}$ . It follows that

$$\begin{array}{ccc} E^{'}D_{TV}^{T}pq^{ptq};p_{i}px^{ptq}qq^{2}/ErReg_{I;T}s/logjIj: \\ t1 & p \end{array}$$

Jensen's inequality now gives that

$$E \overset{T}{\underset{t1}{\overset{}{\stackrel{}{\overset{}}{\overset{}}{\overset{}}}}} D_{TV}pp^{\text{\tiny ptq}}; p_ipx^{\text{\tiny ptq}}qq \; / \; \overset{?}{\overset{?}{\overset{}{\overset{}{\overset{}}{\overset{}}{\overset{}}}}} \overset{g}{\underset{t1}{\overset{}{\overset{}{\overset{}}{\overset{}}}}} D_{TV}pp^{\text{\tiny ptq}}; p_ipx^{\text{\tiny ptq}}qq^2 \; / \; \; \overset{a}{\underset{t1}{\overset{}{\overset{}{\overset{}}{\overset{}}{\overset{}}}}} \overset{a}{\underset{t1}{\overset{}{\overset{}}}} \overset{}{\overset{}{\overset{}}}$$

#### D.2. Proof of Theorem 3.2

Proof of Theorem 3.2. Fix n P N, which we recall represents an upper bound on the description length of the Markov game. Assume that we are given an algorithm B that solves the pT;q-SparseMarkovCCE problem for Markov games G satisfying |G|/n in time U. We proceed to describe an algorithm which solves the 2-player ptn<sup>1{2}</sup>{2u;4q-Nash problem in time pnT U q<sup>Co</sup>, as long as T expp<sup>2</sup> n<sup>1{2}</sup>{2<sup>5</sup>q. First, define n<sub>0</sub>:tn<sup>1{2}</sup>{2u, and consider an arbitrary 2-player n<sub>0</sub>-action normal form G, which is specified by payoff matrices M<sub>1</sub>;M<sub>2</sub> P r0;1s<sup>n<sub>0</sub>n<sub>0</sub></sup>, so that all entries of the game can be written in binary using at most n<sub>0</sub> bits (recall, per footnote 18, that we may assume that the entries of an instance of pn<sub>0</sub>;4q-Nash can be specified with n<sub>0</sub> bits). Based on G, we construct a 2-player Markov game G:GpGq as follows:

Definition D.3. We define the game GpGq to consist of the tuple GpGqpS;H;pA<sub>i</sub>q<sub>iPr2s</sub>;P;pR<sub>i</sub>q<sub>iPr2s</sub>;q, where:

- The horizon of G is  $H 2tn_0 \{ 2u \text{ (i.e., the largest even number at most } n_0 \}$ .
- Let An<sub>0</sub>; the action spaces of the 2 agents are given by A<sub>1</sub> A<sub>2</sub> rAs.
- There are a total of A<sup>2</sup> 1 states: in particular, there is a state s<sub>pa<sub>1</sub>;a<sub>2</sub>q</sub> for each pa<sub>1</sub>;a<sub>2</sub>q PrAs<sup>2</sup>, as well as a distinguished state s, so we have:

$$S tsuYts_{pa^1:a^2q} : pa_1;a_2qPrAs^2u:$$

- For all odd h P rHs, the reward to agents j P r2s given that the action profile pa<sub>1</sub>;a<sub>2</sub>q is played at step h is given by R<sub>j;h</sub>ps;pa<sub>1</sub>;a<sub>2</sub>qq: H Mjqa<sub>1</sub>;a<sub>2</sub>, for all s PS. All agents receive 0 reward at even steps h P rHs.
- At odd steps h P rHs, if actions a<sub>1</sub>;a<sub>2</sub> P rAs are taken, the game transitions to the state s<sub>pa ½</sub>a <sub>2</sub>q. At even steps h P rHs, the game always transitions to the state s.
- The initial state (i.e., at step h1) is s (i.e., is a singleton distribution supported on s).

It is evident that this construction takes polynomial time, and satisfies  $|G|/A^2| 1/n^2|_{0}1/n$ . We will now show by applying the algorithm B to G, we can efficiently compute 4-approximate Nash equilibrium for the original game G. To do so, we appeal to Algorithm 1.

Algorithm 1 Algorithm to compute Nash equilibrium used in proof of Theorem 3.2.

- 1: Input: 2-player, n<sub>0</sub>-action normal form game G.
- 2: Construct the 2-player Markov game GGpGq per Definition D.3, which satisfies |G|/n.
- 3: Call the algorithm B on the game G, which produces a sequence p1q;:::;pTq, where each ptq P markov.
- 4: for tPrTs and odd h PrHs: do
- 5: if ptqpsqPpA1qpA2q is a p4;nq-Nash equilibrium of G: then 6:

return <sup>ptq</sup>psq. <sub>h</sub>

- 7: end if
- 8: end for
- 9: if the for loop terminates without returning: return fail.

Algorithm 1 proceeds as follows. First, it constructs the 2-player Markov game GpGq as defined above, and calls the algorithm B, which returns a sequence  $^{p1q};...;^{pTq}$  P  $^{markov}$  of product Markov policies with the property that the average :  $^{1}$   $^{-}$   $^{T}$   $^{T}$ 

Lemma D.4 (Correctness of Algorithm 1). Consider the normal form game G and the Markov game GGpGq as constructed above, which has horizon H. For any  $_0$  i  $_0$ , TPN, if T exppH $^2$ { $_0^8$ q and  $_1^{p1q}$ ;...; $_1^{p1q}$ P $^{markov}$  are product Markov policies so that  $_1^{t}$  I  $_1^{t}$ I  $_1^{t}$ I  $_2^{t}$ I is an  $_1^{t}$ I  $_2^{t}$ I is an  $_2^{t}$ I  $_2^{t}$ 

The proof of Lemma D.4 is given below. Applying Lemma D.4 with  $_0$  4 (which is a valid application since T exppn $_0$ p4q $^2$ {2 $^8$ q by our assumption on T;), yields that Algorithm 1 always finds a 4-Nash equilibrium of the n $_0$ -action normal form game G, thus solving the given instance of the pn $_0$ ;4q-Nash problem. Furthermore, it is straightforward to see that Algorithm 1 runs in time U pnTq $^{C_0}$ /pUnTq $^{C_0}$ , for some constant C $_0$ ¥1.

Proof of Lemma D.4. Consider a sequence of product Markov policies  $^{p1q}$ ;  $\dots$ ;  $^{p1q}$  with the property that the average  $^{-1}$   $_{T}$   $_{T}$   $_{T}$   $_{T}$   $_{T}$   $_{T}$   $_{T}$   $_{T}$  is an  $p_{0}$ {4q-CCE of G. For all odd h P rHs and j P r2s, let  $p_{j;h}^{p1q}$  psq, P pA<sub>j</sub>q, which is the distribution played under  $^{p1q}$  by player j at step h (at the unique state s with positive probability of being reached at step h). For odd h, we have  $^{p1q}$  psq  $^{p1q}$  p  $^{p1q}$ , and our goal is to show that for some odd h P rHs and t P rTs,  $^{p1q}$  p  $^{p1q}$  is an  $^{0}$ -Nash equilibrium of G. To proceed, suppose for the sake of contradiction that this is not the case.

Let us write  $O_H$ : th P rHs: h oddu to denote the set of odd-numbered steps, and  $E_H$  rHsz $O_H$  to denote the set of even-numbered steps. Let  $H_0$   $|O_H|$   $|E_H|$  H{2. We first note that for j P r2s, agent j's value under the mixture policy is given as follows:

$$V_j^- T \frac{1}{H_{t1}} \frac{1}{h_0} \frac{1}{O} \frac{1}{E_{t1}} \frac{1}{P_{t1}} \frac{1}{P_{t2}} \frac{1}{A_0} \frac{1}{A$$

For each j P r2s, we will derive a contradiction by constructing a (non-Markov) deviation policy for player j in G, denoted  $_j$  P  $_j$ 

Fix any  $h_0$  P rHs,  $j_{;h}$   $\mathbf{1}_0$  P H  $j_{;h}$   $\mathbf{1}_0$  and  $\mathbf{s}_h$  P S. If  $j_{;h}$   $\mathbf{1}$  occurs with positive probability under the transitions of G, then for each h P O<sub>H</sub>, h  $h_0$  1 and both  $j^1$  P r2s, the action played by agent  $j^1$  at step h is determined by  $j_{;h}$ . Namely, if the state at step h 1 of  $j_{;h}$  1 is  $s_{pa^1;a^1q}$ , then player  $j^1$  played action  $a^1$  at step h. So, for each h P O<sub>H</sub> with h  $h_0$  1, we may define  $pa_{1;h}$ ; $a_{2;h}$ q as the action profile played at step h, which is a measurable function of  $j_{;h_0}$ 1. With this in mind, we define  $j_{i;h_0}$ 1; $j_{;h_0}$ 1; $j_{;h_0}$ 1; $j_{;h_0}$ 3 by applying Vovk's aggregating algorithm (Proposition D.2) as follows.

- 1. If h<sub>0</sub> is even, play an arbitrary action (note that the actions at even-numbered steps have no influence on the transitions or rewards).
- 2. If  $h_0$  is odd, define  $\mathbf{p}_{j;h_0} \operatorname{PpA}_j \mathbf{q}$ , by  $\mathbf{p}_{j;h_0} : \operatorname{E}_{\mathsf{t} \mathfrak{q}_{j;h_0}} \operatorname{rp}_{j;h} \mathbf{s}$ , where  $\mathbf{q}_{j;h_0} \operatorname{PprTs} \mathbf{q}$  is defined as follows: for  $\operatorname{tPrTs}$ ,

Note that  $\mathbf{p}_{j;h_0}$  is a function of  $_{j;h_0}$  via the action profiles  $tpa_{1;h};a_{2;h}qu$   $_{h_0:hPO_H}$ ; to simplify notation, we suppress this dependence.

3. Then for any state  $s_{h_0}$  PS, define  $i_{j,h_0}p_{j,h_01}$ ;  $s_{h_0}q$  to be a best response to  $q_{j,h_0}p$  namely

$$: \underset{j;h_0}{p_{j;h}} \underset{1;s_h}{\text{1;sh}} \text{q:argmaxE}_{a} \underset{ajPA_j}{\textbf{p}} \qquad \underset{j;h_0}{rR_{j;h}} \text{ps}_{h_0}; \text{pa1;a2qqsargmaxE}_{a} \underset{ajPA_j}{\textbf{p}} \qquad \underset{j;h_0}{\text{pr}} \text{pM}_{j} \text{qa}_{1}; \text{a2} \text{ss:} \tag{4}$$

Note that, for odd  $h_0$ , the distribution  $pq_{j; h}$   $PpA_{j}q$  defined above can be viewed as an application of Vovk's online aggregation algorithm at step  $ph_0$   $ph_0$  ph

<sup>&</sup>lt;sup>19</sup>Here j denotes the index of the player who is not j.

 $p^{pm}$  phq: $p^{pm}$  Then, the distribution  $q_{j}$ ;p is obtained by updating the aggregation algorithm with the context-observation pairs ph; $a_{j}$ ;h for odd values of h  $h_0$ .

We next analyze the value of  $y_j^{j;j}$  for j Pr2s to show that the deviation strategy we have defined indeed obtains significant gain. To do so, recall that this value represents the payoff for player j under the process in which we draw an index t PrTs uniformly at random, then for each step h PrHs, player j plays according to j and player j plays according to  $p_{j;h}$ .) We recall that  $E_{ij}$  rs denotes the expectation under this process. We let j;h1 PH $_{j;h1}$  denote the random variable which is the history observed by player j in this setup, i.e., when the policy played is j, and let tpa $_{1;h}$ ;a $_{2;h}$ qu denote the action profiles for odd rounds, which are a measurable function of each player's trajectory.

We apply Proposition D.2 with the time horizon as  $H_0$ , and with the set of experts set to  $I:tp^{p1q}; :::p^{p^{T}q}u$  as defined above. The context sequence the sequence of increasing values of  $h P O_H$ , and for each  $h P O_H$ , the outcome at step p0 1q{2 (for which the context is p1 is distributed as p2 in p4 conditioned on p3, which in particular satisfies the realizability assumption stated in Proposition D.2. Then, since (as remarked above), the distributions p4, for p6, are exactly the predictions made by Vovk's aggregating algorithm, Proposition D.2 gives that p2 in p3 in p4 in p5 in p6 in p7 in p8 in p9.

$$E_{:_{j}} = D_{TV}pp_{;h}; p_{j;h}qE_{:_{jpt}q} = D_{TV}pp_{;h}; p^{ptq}phqq^{2}H_{0}logT:$$

$$= D_{TV}pp_{j;h}qE_{:_{jpt}q} = D_{TV}pp_{j;h}; p^{ptq}phq^{2}H_{0}logT:$$

$$= D_{TV}pp_{j;h}qE_{:_{jpt}q} = D_{TV}pp_{j;h}; p^{ptq}phq^{2}H_{0}logT:$$

Recall that we have assumed for the sake of contradiction that  $p_{1;h}^{ptq}p^{ptq}$  is not an 0-Nash equilibrium of G for each h P rHs and t PrTs. Consider a fixed draw of the random variable t PrTs defined above. Then it holds that for j Pr2s and h P rHs, defining

$$o_{;j;h}: \max_{a_1, RrA, pj;_h} rpM_j q_{a_1;a_2} s E^a_{1p1;_h 2p_{2;h}} rpM_j q_{a_1;a_2} s;$$
(6)

$$E_{jp^{j};h}^{a}{}^{pt}p_{s}M_{j}q_{h;j}^{j;h}{}^{j;h}{}^{j;h}{}^{j;h}{}^{j}{}^{j;h}{}^{j}{}^{j}{}^{j}{}^{h}{}^{j}{}^{j}{}^{j}{}^{h}{}^{j}{}^{j}{}^{j}{}^{j}{}^{h}{}^{j}{}^{$$

Combining (6) and (7), we get that for any fixed h PO<sub>H</sub>, j Pr2s, and <sub>i;h1</sub> PH<sub>i;h1</sub>,

$$E^{a}_{j_{p}j_{j_{h}}p_{t}}p_{t}M_{j}q_{j_{j_{h}}p_{j_{j_{h}}}}^{j_{j_{h}}j_{j_{h}}}|t;_{j_{h}1}E^{a}_{j_{p_{1}j_{h}}-2p_{2j_{h}}}rpM_{j}q_{a_{g};a_{2}}si_{\varrho;j_{q}h}2_{j_{j_{h}}}; \tag{8}$$

<sup>20</sup>In fact, Proposition D.2 implies that a similar bound holds uniformly for each possible realization of t, but (5) suffices for our purposes.

Averaging over the draw of t PrTs, which we recall is chosen uniformly, we see that

$$\frac{1}{T} \int_{t_{1}p_{r}^{2}s}^{T} V_{j}^{j} p_{r}^{p_{t}^{q}} V_{j} \qquad p_{t_{q}^{q}} \qquad (9)$$

$$\frac{1}{T}, \frac{T}{t_{1jpr2s}}, \frac{P_{tq}}{r_{t_{1jpr2s}}}, \frac{E}{r_{1j}}, \frac{P_{tq}}{r_{1jh}}, \frac{F_{0j;h}}{r_{1j;h}}, \frac$$

$$\frac{1}{T\,H} \int_{t_1-h_PO_H}^{T,} E_{:ptq} \int_{j}^{r} E_{j,h}^{ptq} rpM_j q_{i-p_{j,h}}^{ptq} rpM_j q_{i-p_{j,h}}^{ptq} rpM_j q_{i-p_{j,h}}^{ptq} rpM_j q_{a_1;a_2} s_j pr2s_{a_1,b_1}^{ptq} rpM_j q_{a_1;a_2} s_j pr2s_{a_1,b_2}^{ptq} rpM_j q_{a_1;a_2}^{ptq} rpM_j q_{a_1;a$$

$$\frac{1}{T H} \int_{t_{1}Pr2s}^{T} \int_{j=j-hPO_{H}}^{T} \int_{t_{1}Pr2s}^{T} \int_{j}^{t_{1}Pr2} \int_{j}$$

$$\underbrace{ \stackrel{0}{2}}_{\mathsf{T}} \underbrace{ \stackrel{1}{\mathsf{H}}_{\mathsf{+}\mathsf{1}}}^{2} \stackrel{1}{\mathsf{H}_{\mathsf{0}}} \underbrace{ \underset{0}{\mathsf{logT}}}_{\mathsf{T}} \underbrace{ \stackrel{3}{\mathsf{H}_{\mathsf{0}}} \mathsf{ogpT}}_{\mathsf{q}\{\mathsf{H};}$$
 (11)

where (9) follows from the definition  $\frac{1}{2}$  T  $I_{\frac{1}{2}}$ , (10) follows from (8), and (11) uses (5). As long as T exppH  $p_0$ {16 $q^2q$ , the this expression is bounded below by  $_0$ {4, meaning that is-not an  $_0$ {4-approximate CCE. This completes the contradiction.

#### E. Proofs of lower bounds for Sp<sup>A</sup> RSECCE (Sections 4 and 5)

In this section we prove our computational lower bounds for solving the SPARSECCE problem with m 3 players (Theorem 4.3 and Corollary 4.4), as well as our statistical lower bound for solving the SPARSECCE problem with a general number m of players (Theorem 5.2).

Both theorems are proven as consequences of a more general result given in Theorem E.1 below, which reduces the NASH problem in m-player normal-form games to the Sparsecce problem in pm 1q-player Markov games. In more detail, the theorem shows that (a) if an algorithm for Sparsecce makes few calls to a generative model oracle, then we get an algorithm for the NASH problem with few calls to a payoff oracle (see Section C.3 for background on the payoff oracle for the NASH problem), and (b) if the algorithm for Sparsecce is computationally efficient, then so is the algorithm for the NASH problem.

Theorem E.1. There is a constant  $C_0 \neq 0$  so that the following holds. Consider n;m P N, and suppose T;N;Q P N and  $\neq 0$  satisfy 1 T  $\exp \frac{2 \operatorname{tn} \{ mu \}}{m^2}$ . Suppose there is an algorithm B which, given a generative model oracle for a pm 1q-player Markov game G with |G|/n, solves the pT;;N q-SPARSECCE problem for G using Q generative model oracle queries. Then the following conclusions hold:

- For any i 0, the m-player ptn{mu;16pm 1qq-Nash problem for any normal-form game G can be solved, with failure probability, using at most CopQlogp1{qq plogp1{qnm{q^{Co} queries to a payoff oracle O\_G for G.}
- If the algorithm B additionally runs in time U for some U PN, then the algorithm solving NASH from the previous bullet point runs in time pnmTNUlogp1{q $q^{C_0}$ .

Theorem 4.3 follows directly from Theorem E.1 by taking m2.

Proof of Theorem 4.3. Suppose there is an algorithm which, given the description of any 3-player Markov game G with |G|/n, solves the pT;;Nq-SPARSECCE problem in time U. Such an algorithm immediately yields an algorithm which can solve the pT;;Nq-SPARSECCE problem in time U  $|G|^{Op1q}$  using only a generative model oracle, since the exact description of the Markov game can be obtained with |G|/n using only a generative model oracle, since the exact description of the Markov game can be obtained with |G|/n using only a generative model oracle, since the exact description of the Markov game can be obtained with |G|/n queries to the generative model (across all ph;s;aq tuples). We can now solve the problem of computing a 50-Nash equilibrium of a given 2-player tn{2u-action normal form game G as follows. We simply apply the algorithm of Theorem E.1 with m2, noting that the oracle  $O_G$  in the theorem statement can be implemented by reading the corresponding bits of input of the input game G. The second bullet point yields that this algorithm

expp<sup>2</sup>tn{mu{m<sup>2</sup>q of Theorem E.1 is takes time pnT NU logp1{q{q<sup> $C_0$ </sup>, for some constant  $C_0$ . Furthermore, the assumption T implied by the assumption that T  $\exp^2 n\{16a \text{ of Theorem 4.3.}$ 

In a similar manner, Theorem 5.2 follows from Theorem E.1 by applying Theorem C.3, which states that there is no randomized algorithm that finds approximate Nash equilibria of m-player, 2-action normal form games in time 20pmq.

Proof of Theorem 5.2. Let 0 be the constant from Theorem C.3, and consider any m ¥ 3. Suppose there is an algorithm which, for any m-player Markov game G with  $|G|/2m^6$ , makes Q oracle queries to a generative model oracle for G, and solves the pT;<sub>0</sub>{p10mq;Nq-SparseCCE problem for G for some T;N PN so that T exppcmq, for a sufficiently small absolute constant c. Then, by Theorem E.1 with o{p10mq and n m<sup>6</sup> (which ensures that T expppo{p10mqq<sup>2</sup> tn{mu{m<sup>2</sup>q as long as c is sufficiently small), there is an algorithm which solves the pm<sup>5</sup>;<sub>0</sub>q-Nash problem—and thus the p2;<sub>0</sub>q-Nash problem—for pm1qplayer games with failure probability 1{3, using OpQq m<sup>Op1q</sup> queries to a payoff oracle. But by Theorem C.3, any such algorithm requires pmq follows that queries oracle. payoff pmq, as desired.

#### E.1. Proof of Theorem E.1

Proof of Theorem E.1. Fix any m ¥ 2, n P N. Suppose we are given an algorithm B that solves the pm pT;;Nq-SparseCCE problem for Markov games G satisfying |G|/ n, running in time U and using at most Q generative model queries. We proceed to describe an algorithm which solves the m-player ptn{mu;16pm 1qq-NASH problem using Co pQlogp1{qq plogp1{q nm{ $q^{C_0}}$  queries to a payoff oracle, and running in time pnmTNU logp1{ $q(q^{C_0})$ , where represents the failure probability. Define no :tn{mu, and assume we are given an arbitrary m-player no-action normal form G, which is specified by payoff matrices  $M_1$ ;:::; $M_m Pr0$ ; $1s^{n_0n_0}$ . We assume that all entries of each of the matrices  $M_i$  have only the most significant maxtno; rlog 1 (su bits nonzero; this assumption is without loss of generality, since by truncating the utilities to satisfy this assumption, we change all payoffs by at most , which degrades the quality of any approximate equilibrium by at most 2 (in addition, we have  $rlog1{s/n_0}$  since we have assumed 1 T  $\exp^2 n_0 \{m^2 q\}$ . We assume  $1 \{2 \text{ without } \}$ loss of generality. Based on G, we construct an pm 1q-player Markov game G:GpGq as follows.

Definition E.2. We define the Markov game GpGq as the tuple GpGq pS;H;pA<sub>i</sub>q<sub>iPr2s</sub>;P;pR<sub>i</sub>q<sub>iPr2s</sub>;q, where:

- The horizon of G is chosen to be the power of 2 satisfying  $n_0/H$
- Let A:n<sub>0</sub>. The action spaces of agents 1;2;:::;m are given by A<sub>1</sub> A<sub>m</sub> rAs. The action space of agent m 1 is

$$A_{\,m} \quad {}_1tpj; a_jq: j\, P\, rms; a_j\, P\, A_j\, u;$$

so that  $|A_{m-1}|Am/n$ .

We write A  $\frac{1}{m}$  A<sub>j</sub> to denote the joint action space of the first m agents, and  $\overline{A}$ :  $\frac{m}{j1}$  A<sub>j</sub> to denote the joint action space of all agents.

- There is a single state, denoted by s, i.e., S tsu (in particular, is a singleton distribution supported on s).
- For all h P r H s, the reward for agent j P r m 1s, given an action profile a pa<sub>1</sub>;:::;  $a_{m-1}q$  at the unique state s, is as follows: writing  $a_m _1 pj^1; a^1_1q$ , we have

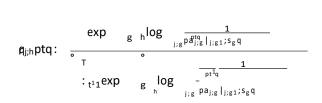
$$R_{j;h}$$
ps;aq $R_{j;h}$ ps;aq  $\frac{1}{2}$ 3rlog1{s} encpaq; (12)

where R<sub>i;h</sub>ps;aq is defined per the kibitzer construction of (Borgs et al., 2008):

In (12) above, encpaqPr0;1s is the binary representation of a binary encoding of the action profile a. In particular, if the binary encoding of a is  $pb_1;...;b_Nq$ , with  $b_i$  Pt0;1u, then encpaq  $\sum_{i=1}^{N} 2^i b_i$ . Note that encpaq takes N Opmlogn<sub>0</sub>q/Opmlognq bits to specify.

Algorithm 2 Algorithm to compute Nash equilibrium used in proof of Theorem E.1.

- 1: Input:
- 2: Parameters  $n;n_0;m;T PN$ , {p6Hq, Kr4logpmn<sub>0</sub>{q{ $}^2$ s.
- 3: An m-player,  $n_0$ -action normal form game G, with utilies accessible by oracle  $O_G$ .
- 4: An algorithm B for computing approximate CCE of Markov games.
- 5: Call the algorithm B on the pm 1q-player Markov game G GpGq constructed as in Definition E.2, which produces a sequence p1q;:::;pTq, where each ptq potq;:::;pTq 1q with ptq p gen;rnd. Here, we use the oracle O<sub>G</sub> to simulate generative model oracle queries made by B.
- 6: Draw t PrTs uniformly at random.
- 7: For each j P rms, initialize  $_{j;0}$  to be an empty trajectory.
- 8: for h P r H s: do
- 9: Set  $s_h$  s (per the transitions of G).
- 10: For each j Prms, define  $q_{j;h}: E_{tq_{j;h}=j;h}p_{j;h}q^{tq}s_hq$  PpA<sub>j</sub>q, where  $q_{j;h}$  PprTsq is defined as follows: for tPrTs,



- 11: Draw K i.i.d. samples  $a_h^1$ ;...; $a_h^{K}$   $j_{Prms}q_{j;h}$ .p
- 12: For each  $a^1 P A_m$  1, define  $R_m$  1, hpa $^1q$ :  $K_1 R_m$  1, hpsh; pa $^k$ ;  $R_m$  1,
- 13: For each j P rms, draw  $a_{j;h_{j;h}}$   $p_{j;h_1;s_hq}$ .
- 14: Choose the action a<sub>m 1;h</sub> of player m 1 as follows: (Action a<sub>m 1;h</sub> is corresponds to the action selected by the policy m<sup>2</sup> of player m 1 defined within the proof of Lemma E.3; this policy is well-defined because the action profiles of all players i Prms can be extracted from the lower-order bits of player m 1's reward)

$$\begin{array}{ccc} & ! & ) \\ a_{m} \ _{1;h} : argmax \ R_{n}p \ _{1;h}pa^{1}q \ : & \\ & & \\ a^{1}PA_{m} \ _{1} \end{array} \tag{14}$$

- 15: For each j P rm 1s, let  $r_{i;h}$   $R_{i;h}$   $ps_h$ ;  $pa_{1;h}$ ; ...;  $a_{m-1;h}$  qq.
- 16: Each player j constructs  $_{j;h}$  by updating  $_{j;h1}$  with ps $_h;a_{j;h};r_{j;h}q$ . 17: if  $R_{m-1;h}pa_{m-1;h}q/14pm$  1q{H then
- 18: return  $\mathbf{p}_h$ :  $\mathbf{p}_{j;h}$  as a candidate approximate Nash equilibrium for G.
- 19: end if
- 20: end for
- 21: if the for loop terminates without returning: return fail.

It is evident that this construction takes polynomial time and satisfies  $|G|/mn_0/n$ . Furthermore, it is clear that a single generative model oracle call for the Markov game G (per Definition 5.1) can be implemented using at most 2 calls to the oracle  $O_G$  for the normal-form game G. We will now show by applying the algorithm B to G, we can efficiently (in terms of runtime and oracle calls) compute a 16pm 1q-approximate Nash equilibrium for the original game G. To do so, we appeal to Algorithm 2.

Algorithm 2 proceeds as follows. First, it calls the algorithm B on the pm 1q-player Markov game GpGq, using the oracle  $O_G$  to simulate B 's calls to the generative model oracle for G. By assumption, the algorithm B returns a sequence  $^{p1q}$ ;:::; $^{pTq}$  of product policies of the form  $^{ptq}$   $p^{ptq}$ ;:::; $^{pTq}$   $_1$ q,  $_1$ q,  $_2$ q,  $_3$ 0 that each  $^{ptq}$   $P^{gen;rnd}$  is  $_1$ -computable, and so that the average :  $_1$ 

is an -CCE of G. Next, Algorithm 2 samples a trajectory from G in which:

- Players 1;:::;m each play according to a policy ptq for an index t PrTs chosen uniformly at the start of the episode.
- Player m 1 plays according to a strategy that, at each step h P rHs, computes distributions  $q_{j;h}$  representing its "belief" of what action each player j P rms will play at step h (Line 10), and plays an approximate best response to the product of the strategies  $p_{i;h}$ , j P rms (Line 14).

In order avoid exponential dependence on the number of players  $m_i$  when computing an approximate best response to  $^{j}_{Prms}q_{j;h}p$  we draw K:r4logpmn<sub>0</sub>{q{ $^2$ s (for {p6Hq}) samples from  $^{i}_{Prms}g_{j;h}$  and use these samples to compute the best response.

In particular, letting  $a_h^K PA$  denote the kth sampled action profile, we construct a function  $R_{a=1;h}: A_{m=1}\tilde{N}$  R in Lines 11 and 14 which, for each  $a^1 PA_{m=1}$ , is defined as the average over samples  $ta^k u_{kR_frKs}$  of the realized payoffs  $R_{m=1;h}ps_h;pa^k;a^1qq;$  note that to compute the payoffs for each sample, Algorithm 2 needs only two oracle calls to  $O_G$ .

The following lemma, proven in the sequel, gives a correctness guarantee for Algorithm 2.

Lemma E.3 (Correctness of Algorithm 2). Given any m-player  $n_0$ -action normal form game G, if the algorithm B solves the pT;;Nq-SparseCCE problem for the game GpGq with T;;N satisfying T/exppn $_0^2$ {m $^2$ q, then Algorithm 2 outputs a 16pm 1q-approximate Nash equilibrium of G with probability at least 1{3, and otherwise fails.

The assumption that  $T = \exp^{2tn\{mu\}}$  from the statement of Theorem E.1 yields that  $T/\exp pn_0^2\{m^2q$ , so Lemma E.3 yields that Algorithm 2 outputs a  $16pm^{m^2}1q$ -Nash equilibrium of G with probability at least 1{3 (and otherwise fails). By iterating Algorithm 2 for logp1{q times, we may thus compute a 16pm 1q-Nash equilibrium of G with failure probability 1.

We now analyze the oracle cost and computational cost of Algorithm 2. It takes 2Q oracle calls to  $O_G$  to simulate the Q generative model oracle calls of B , and therefore, if B runs in time U, then the call to B on Line 5, using oracle calls to  $O_G$  to simulate simulate the generative model oracle calls, runs in time OpUq. Next, the computations of  $q_{jrh}$  (and thus  $q_{j;h}$ ) in Line 10 can be performed in pnmTNq<sup>Op1q</sup> time, the computation of  $R_{TD}$  1;h:  $A_{m-1}$   $\tilde{N}$  R in Line 14 requires time (and oracle calls to  $O_G$ ) bounded above by Op| $A_{m-1}$ |Kq/pnmlogp1{q{ $q^{Op1q}$ , constructing the actions  $a_{j;h}$  (for j Prm 1s) in Lines 13 and 14 takes time pNmnq<sup>Op1q</sup> (using the fact that the policies  $a_{j;h}$  are  $a_{j;h}$  are  $a_{j;h}$  on Line 15 requires another 2pm 1q oracle calls to  $a_{j;h}$  Altogether, Algorithm 2 requires 2Q pnmlogp1{q{ $a_{j;h}$  oracle calls to  $a_{j;h}$  or Line 15 runs in time U, then Algorithm 2 takes time pnmTNUlogp1{q{ $a_{j;h}$  for some absolute constant  $a_{j;h}$  oracle calls to  $a_{j;h}$  oracle calls to  $a_{j;h}$  or some absolute constant  $a_{j;h}$  oracle calls to  $a_{j;h}$  oracle calls to  $a_{j;h}$  or some absolute constant  $a_{j;h}$  oracle calls to  $a_{j;h}$  oracle calls to

Remark E.4 (Bit complexity of exponential weights updates). In the above proof we have noted that  $q_{j;h}$  (as defined in Line 10 qf] Algorithm 2) can be computed in time pnmTNq<sup>Op1q</sup>. A detail we do not handle formally is that, since the values of  $q_{j;h}$ ptq are in general irrational, only the pnmTNq<sup>Op1q</sup> most significant bits of each real number  $q_{j;h}$ ptq can be computed in time pnmTNq<sup>Op1q</sup>. To give a truly polynomial-time implementation of Algorithm 2, one can compute only the pnmTNq<sup>Op1q</sup> most significant bits of each distribution  $q_{j;h}$ , which is sufficient to approximate the true value of  $q_{j;h}$  to within expppnmTNq<sup>Op1q</sup>  $q_{j;h}$  to total variation distance. Since  $q_{j;h}$  only influences the subsequent execution of Algorithm 2 via the samples  $q_{j;h}$  in total variation distance. In particular, the correctness guarantee of Lemma E.3 still holds, with sucess probability at least 1{3 expppnmTNq<sup>Op1q</sup>q i 1{4}.

It remains to prove Lemma E.3, which is the bulk of the proof of Theorem E.1.

Proof of Lemma E.3. We will establish the following two facts:

- 1. First, the choices of a<sub>m 1;h</sub> in Line 14 (i.e., Eq. 14) of Algorithm 2 correspond to a valid policy important for player m 1 (representing a strategy for deviating from the equilibrium ), in that they can be expressed as a function of player pm 1q's history, p<sub>m 1;h1</sub>; s<sub>h</sub>q at each step h.
- Second, we will show that, since is an -CCE of G, the strategy in a cannot not lead to a large increase of value for player must return a Nash equilibrium with high enough probability.

Defining  $\frac{1}{2}$  for i P rm 1s. We begin by constructing the policy  $\frac{1}{2}$  described; for later use in the proof, it will be convenient to construct a collection of closely related policies Pen;rnd for iPrms, also representing strategies for deviating from the equilibrium .-

Let i P rm  $\,$  1s be fixed. For h P rHs, the mapping  $\,^{:}_{\,\,i;h}:H_{i;h1}\,$  S  $\,\tilde{N}\,$  A $_{i}$  is defined as follows. Given a history  $_{i;h1}$  $ps_1; a_{i;1}; r_{i;1}; ...; s_{h1}; a_{i;h1}; r_{i;h1} \neq PH_{i;h1}$  (we assume without loss of generality that i;h1 occurs with positive probability under some sequence of general policies) and a current state  $s_h$ , we define  $_{i:h}p_{i;h1}$ ;  $s_hqPA_i$  through the following process.

- 1. First, we claim that for all players j Prm 1sztiu, it is possible to extract the trajectory i;h1 from the trajectory i;h1 of player i.
  - (a) Recall that for each g h, from the definition in (12) and the function encpaq, the bits following position 3rlog1{s of the reward  $r_{i;g}$  given to player i at step g of the trajectory  $_{i;g1}$  encode an action profile  $a_g P A$ . Since  $_{i;h1}$  occurs with positive probability, this is precisely the action profile which was played by agents at step g. Note we also use here that by definition of the rewards  $R_{i;h}$ ps;aq in (12), the component  $R_{i;h}$ ps;aq  $\overline{of}$  the reward only affects the first 2rlog1{s bits.
  - (b) For g h and j Prm 1sztiu, define  $r_{j;g}:R_{j;g}ps_g;a_gq$ .
  - (c) For j Prm 1sztiu, write j;h1:ps1;aj;1;rj;1;:::;sh1;aj;h1;rj;h1q; in particular, j;h1 is a deterministic function of pi;h1;shq. (Note that, since i;h1 occurs with positive probability, the history i;h1 observed by player jup to step h1 can be computed from it via Steps (a) and (b)). Going forward, for g = h1, we let  $_{j;g}$  denote the prefix of  $_{j;h1}$  up to step g.
- 2. Now, using that player i can compute all players' trajectories, for each j Prm 1s we define

$$\mathbf{p}_{j;h}: \mathbf{E}_{\mathsf{t}\mathfrak{q}_{j;h}} :_{i:h} \mathbf{p}_{j;h}^{\mathsf{p}_{i}\mathsf{q}}, \mathbf{s}_{h} \mathbf{q} \quad \mathsf{P} \mathsf{p} \mathsf{A}_{j} \mathbf{q}; \tag{15}$$

where q<sub>i;h</sub> PprTsq is defined as follows: for tPrTs,

paj;g | j;g1;sgq

Note that  $\mathbf{p}_{j;h}$  is a random variable which depends on the trajectory  $\mathbf{p}_{j;h1}$ ;  $\mathbf{s}_{hq}$  (which can be computed from  $\mathbf{p}_{i;h1}$ ;  $\mathbf{s}_{hq}$ ). In addition, the definition of  $\mathbf{p}_{jh}$  (for each j Prms) is exactly as is defined in Line 10 of Algorithm 2.

3. For iPrms, define  $_{i;h}p_{i;h1}$ ; shq as follows:

For the case i m 1, define  $a_{m-1;h} p_{m-1;h1}$ ;  $s_h q P p A_{m-1} q$  (implicitly) to be the following distribution over  $a_{m-1;h} P \stackrel{:}{A}_{m-1}$ : draw  $a^1; :::; a^K_h$   $_h$   $_{jPrms} p_{j;h}$ , define  $p_{m-1;h} p_{m-1;h} p_{m-1;h}$ 

$$a_{m-1;h}^{:}$$
:  $a_{n}^{:}$   $a_{n}^{:}$  (18)

Note that, for each choice of  $p_{m-1;h1}$ ;  $s_hq$ , the distribution  $\frac{1}{m-1;h1}$ ,  $p_{m-1;h1}$ ;  $s_hq$  as defined above coincides with the distribution of the action  $a_{m-1;h}^{:}$  defined in Eq. 14 in Algorithm 2, when player m 1's history is  $m_{1;h}$  and the state at step h is , $s_h$ . The following lemma, for use later in the proof, bounds the approximation error incurred in sampling aդ;:::;aқ jPrms ₽j;h.

 $\text{Lemma E.5. Fix any } p_{m-1;h1}; s_h q \ P \ H_{j;h1}. \ \text{With probability at least 1 over the draw of } a^1; \dots; a^K_{h} \stackrel{\cdot}{j}_{Prms}, q_{j;h},$ р it holds that for all a<sup>1</sup>P A<sub>m</sub> <sub>1</sub>,

$$RP_{m-1;h}pa^1qE_{a_j\phi^i_{jh}} = rR_{m-1;h}ps_h;pa_1;:::;a_m;a^1qqs_H$$
; —

which implies in particular that with probability at least 1 over the draw of a \_\_\_\_\_1;h m\_\_\_1;h p\_m\_\_\_1;h1;shq,

$$\max_{a^1p_{A_m}} E_{a_j \phi_{j;h} @ jPrms} rR_{m-1;h} ps_h; pa_1; ...; a_m; a^1qqs (\frac{2}{H} E_{a_j \phi_{j;h} @ jPrms} rR_{m-1;h} ps_h; pa_1; ...; a_m; a_{m-1}^{'}{}_{h} \dot{q}qs:$$
 (19)

It is immediate from our construction above that the following fact holds.

Lemma E.6. The joint distribution of j;h, for j Prm 1s and h PrHs, as computed by Algorithm 2, coincides with the distribution of j;h in an episode of G when players follow the policy  $m = \frac{1}{1}pm = \frac{1}{10}$ .

Analyzing the distributions  $\mathbf{q}_{j;h}$ . Fix any i P rm 1s. We next prove some facts about the distributions  $\mathbf{q}_{j;h}$  defined above (as a function of  $p_{i;h1}$ ; $s_hq$ ) in the process of computing  $s_{i;h}p_{i;h1}$ ; $s_hq$ .

For each h P r Hs, consider any choice of  $p_{i;h1}$ ;  $s_hq$  P  $H_{i;h1}$ S; note that for each j P r m 1s, the distributions  $q_{j;h}$  P  $p_{Aj}q$  for h P r Hs may be viewed as an application Vovk's aggregating algorithm (Proposition D.2) in the following setting: the number of steps (T, in the context of Proposition D.2; note that T has a different meaning in the present proof) horizon is H, the context space is  $h_1^H H_{j;h1} S$ , and the output space is  $A_j$ . The expert set is I  $t^{p_1q}$ ;...; $p^{r_1q}u_j$  (which has jIj T), and the experts' predictions on a context  $p_{j;h1}$ ;  $s_jq$  P  $H_{j;h1}S$  are defined via  $p_{j;h1}$ ;  $s_jq$  P  $p_{j;h1}$  P  $p_{j$ 

In more detail, fix any tPrTs and jPrm 1s with ij. We may apply Proposition D.2 with the number of steps set to H, the set of experts as I<sub>j</sub> t<sup>p1q</sup>;:::;<sup>pTq</sup>u, and contexts and outcomes generated according to the distribution induced by running the policy i iq in the Markov game G as follows:

- For each h PrHs, we are given, at steps  $h^1$  h, the actions  $a_{k;h^1}$  rewards  $r_{k;h^1}$  for all agents k Prm 1s, as well as the states  $s_1;...;s_h$ .
  - For each k Prm 1s, set  $_{k;h1}$  ps<sub>1</sub>; $a_{k;1}$ ; $r_{k;1}$ ;...; $s_{h1}$ ; $a_{k;h1}$ ; $r_{k;h1}$ q to be agent k's history. The context fed to the aggregation algorithm at step h is  $p_{i:h1}$ ; $s_{hq}$ .
  - The outcome at step h is given by  $a_{j;h}^{ptq}p|_{\substack{j;h\\j;h}}$ ;  $s_hq$ ; note that this choice satisfies the realizability assumption in Proposition D.2.
  - To aid in generating the next context at step h 1, choose ak;h pkth;shq for all kPrm 1szti;ju and ai;h
     i;hpi;h1;shq. Then set sh 1 to be the next state given the transitions of G and the action profile ah pa1;h;:::;am 1;hq.

By Proposition D.2, it follows that for any fixed tPrTs and jPrm 1s with ji, under the process described above we have

Analyzing the value of  $_{m}^{:}$  1. Next, using the development above, we show that if Algorithm 2 successfully computes a Nash equilibrium with constant probability (via  $_{m}^{:}$  1) whenever s is an -CCE. We first state the following claim, which is proven in the sequel by analyzing the values  $V_{i}^{:}$  for iPrms.

Lemma E.7. If is an -CCE of G, then it holds that for all iPrms,

$$V_i^ = \frac{a}{\log pT} \frac{a}{q\{H:}$$

Note that in the game G, since for all h P rHs, s P S and  $\overline{aP}$  A, it holds that  $m_{j1}^{n} R_{j;h} ps; aq / \frac{pm - 1q^2}{H}$  (which holds since in (12), encpaq is multiplied by  $m_{j1} 2^{3r\log 1\{s\}}$ ), it follows that  $m_{j1} 2^{-1} V / pm - 1q^2$ . Thus, by Lemma E.7, we have  $m_{j1} 2^{-1} pm - 1q^2 pm - 1q^2$ 

$$V_{m_1}^{:p^{-1}p^{-1}q}/2pm \quad 1q \quad m^2 \frac{a}{\log pT} q \{H:$$
 (21)

To simplify notation, we will write  $\mathbf{q}_h : \mathbf{p}_{;h} \mathbf{q}_{m;h}$  in the below calculations, where we recall that each  $\mathbf{q}_{j;h}$  is determined given the history up to step h,  $\mathbf{p}_{j;h1}; \mathbf{s}_h \mathbf{q}$ , as defined in (15) and (16). An action profile drawn from  $\mathbf{q}_h$  is denoted as a  $\mathbf{q}_h$ , with a P.A. We may now write V  $\mathbf{q}_h^{1^{pm-1}\mathbf{q}}$  as follows:

$$V_{\text{trTs}}^{\text{mi}} = \frac{1}{H} \frac{1}{H} = \frac{1}{H} \frac{1}$$

where:

- The first inequality follows from the fact that R<sub>m</sub> 1;hpq takes values in r1{H;1{Hs and the fact that the total variation between product distributions is bounded above by the sum of total variation distances between each of the pairs of component distributions.
- The second inequality follows from the inequality (19) of Lemma E.5.
- The final equality follows from the definition of the rewards in (12) and (13), and by summing (20) over j P rms. We remark that the <sup>2</sup> term in the final line comes from the term H 2<sup>3</sup>/<sub>4</sub> rlog1{s encpaq in (12).

Rearranging and using (21) as well as the fact that  $^2\{p6Hq\ ^2/\ (as/1\{2),\ we\ get\ that\ ^2\}\}$ 

$$E_{t rTs} E_{i pt q}, max E_{ap h} pM_{j} q_{a_{j};a_{j}} pM q_{s} s_{a_{j}}$$

$$= \frac{1}{2} \frac{1$$

Since  $\mathbf{p}_h$  is a product distribution a.s., we have that

$$\max_{j\, P\, r\, m\, s;\, a_{i_i,h}\, P\, A_j} \mathsf{E}_{a_{\boldsymbol{p}} \quad _h} \, r\! p\! M_j \, q_{a_j^{\, 1};a_j} \quad p\! M \, \, q_j \, s_a^{\, 2} \, 0 \colon$$

Therefore, by Markov's inequality, with probability at least 1{2 over the choice of t rTs and the trajectories  $p_{j;h1}$ ;  $s_h q_{m-1} \frac{1}{pm-1q}$  for J'P'rms (which collectively determine  $q_h$ ), there is so that

max 
$$E_{ap,h} rpM_j q_{j;a_j}^1$$
  $pM q_s/10pm$  1q 2pm 1qm  $\frac{a}{logp} Tq{H/12pm}$  1q; (22)

where the final inequality follows as long as H  $^2$  ¥ m $^2$ logT, i.e., T / exp  $^H$   $\frac{^2}{m^2}$  , which holds since H ¥ n $_0$  and we have assumed that T / expp $^2$  n $_0$ {m $^2$ q.

Note that (22) implies that with probability at least 1{2 under an episode drawn from  $\frac{1}{m-1}$  pm  $\frac{1}{1}$ q, there is some h P r Hs so that qh js a 12pm  $\frac{1}{1}$ q-Nash equilibrium of the stage game G. Thus, by Lemma E.6, with probability at least 1{2 under an episode drawn from the distribution of Algorithm 2, there is some h P r Hs so that  $\mathbf{p}$  is a 12pm  $\frac{1}{1}$ q-Nash equilibrium of G.

Finally, the following two observations conclude the proof of Lemma E.3.

• If  $g_h$  is a 12pm 1q-Nash equilibrium of G, then by definition of the reward function R<sub>m</sub> 1;hpq in (12), upper bounding H  $\frac{2^{3\text{rlog1}(s)}}{5^{3\text{rlog1}(s)}}$  encpag by  $^2$  {H,

$$\max_{a^1} E_{a \phi_h} R_{m 1;h} ps; pa; a^1 qq/H \xrightarrow{12pm} 1q H; \xrightarrow{2}$$

which implies, by Lemma E.5, that with probability at least 1 over the draw of a ;:::; $^1\!a_h$  ,  $^K$ 

! ) 
$$\max_{a^1PA_{m-1}} R_{m-1;h}pa^1q / \frac{1}{H} 12pm 1q + \frac{2}{H} + \frac{1}{H} \frac{1}{44pm} 1q;$$

i.e., the check in Line 17 of Algorithm 2 will pass and the algorithm will return  $\mathbf{p}_h$  (if step h is reached).

• Conversely, if  $\max_{a^1PA_{m-1}} \Re_{m-1/n}pa^1q / 14pm$  1q, i.e., the check in Line 17 passes, then by Lemma E.5, with probability at least 1 over  $a_h$ ;:: $^1_{a_h}$ ,  $^K$ 

$$\max_{a^1} E_{a \cdot \phi_h} R_{m \ 1;h} ps; pa; a^1 qq/H \xrightarrow{14pm} 1q \ H/H \xrightarrow{15pm} 1q;$$

which implies, by the definition of  $R_{m-1;h}$ pq in (12) and (13), that  $p_h$  is a 16pm 1q-Nash equilibrium of G.

Taking a union bound over all H of the probability-failure events from Lemma E.5 for the sampling  $a^1$ ;::;;  $a^K p_h$  (for  $h_h PrHs$ ), as well as over the probability-1{2 event that there is no  $q_h$  which p is a 12pm 1q-Nash equilibrium of G, we obtain that with probability at least 11{2H{p6Hq} \$\frac{1}{4}\$} 13, Algorithm 2 outputs a 16pm 1q-Nash equilibrium of G.

Finally, we prove the remaining claims stated without proof above.

Proof of Lemma E.5. Since  $R_{m-1;h}ps;aq Pr1\{H;1\{Hs for each a PA, \overline{b}y Hoeffding's inequality, for any fixed <math>a^1PA_{m-1}$ , with probability at least  $1\{|A_{m-1}|1\{pmn_0q \text{ over the draw of a };...;a_h^{-1} K_{jPrms}, p_{j;h}, \text{ it holds that}\}$ 

$$RP_{m-1;h}pa^{1}qE_{a_{j}\mathbf{p}_{,h}^{i}}@_{j}Prms}rR_{m-1;h}ps_{h};pa_{1};...;a_{m};a^{1}qqs/H$$

where the final inequality follows from the choice of K r4logpmn<sub>0</sub>{ $q\{^2s\}$ . The statement of the lemma follows by a union bound over all  $|A_m|_1$  actions  $a^1PA_m|_1$ .

Proof of Lemma E.7. Fix any agent i P rms. We will argue that the policy  $_{i}$   $\stackrel{p}{\mapsto}_{i}^{\text{gen;d}}$   $_{i}^{\text{d}}$  defined within the proof of Lemma E.3 satisfies  $V_{i}$   $\stackrel{i}{\mapsto}_{i}$   $\stackrel{a}{\neq}$   $_{i}$   $\stackrel{a}{\mapsto}_{i}$   $\stackrel{a}{\mapsto}_{i}$   $_{i}$   $\stackrel{a}{\mapsto}_{i}$   $\stackrel{a}{\mapsto}_{i}$ 

from which the result of Lemma E.7 follows after rearranging terms. To simplify notation, let us write  $\mathbf{p}_{i;h}$ :  $\mathbf{q}_{i;h}\mathbf{p}$  where we recall that eagh  $\mathbf{q}_{i;h}$  is determined given the history up to step h,  $\mathbf{p}_{i;h1}$ ;  $\mathbf{s}_h\mathbf{q}$ , as defined in (15) and (16). An action profile

drawn from  $\mathbf{p}_{i;h}$  is denoted by  $\mathbf{a}_i \mathbf{p}_{i;h}$ , with  $\mathbf{a}_i P A_i$ . We compute

#### where:

- The first inequality follows from the fact that the rewards R<sub>i;h</sub>pq take values in r1{H;1{Hs and that the total variation between product distributions is bounded above by the sum of total variation distances between each of the pairs of component distributions.
- The second inequality follows from the definition of inplication in plants: pinks in terms of qi;h in (17) as well as (20) applied to each ji and each tPrTs.
- The final inequality follows by Lemma E.8 below, applied to agent i and to the distribution  $\mathbf{p}_{i;h}$ , which we recall is a product distribution almost surely.

 $\label{eq:emma_energy} \text{Lemma E.8. For any iPrms, sPS;hPrHs, and any product distribution qPpA$_{i}\overline{q}$, it holds that $\max E_{aq}R_{i;h}ps;pa_{i};aqq$0:$_{i}$$ 

Proof. Choose  $a_i$ :  $argmax_{a^1PA}E_{aq}pM_iq_{a^1;a}$ . Now we compute

where the first inequality follows since q is a product distribution, the second inequality uses that encpq is non-negative, and the final inequality follows since by choice of  $a_i$  we have  $E_{aq}$   $pM_iqa_i$ ;  $a_i \neq E_{aq}$   $pM_iqa_{a;a}$  for all  $a_i \neq A_i$ .

#### E.2. Remarks on bit complexity of the rewards

The Markov game GpGq constructed to prove Theorem E.1 uses lower-order bits of the rewards to record the action profile taken each step. These lower order bits may be used by each agent to infer what actions were taken by other agents at the previous step, and we use this idea to construct the best-response policies idefined in the proof. As a result of this aspect of the construction, the rewards of the game GpGq each take Opmlogpnq logp1{qq bits to specify. As discussed in the proof of Theorem E.1, it is without loss of generality to assume that the payoffs of the given normal-form game G take Oplog1{q bits each to specify, so when either m " 1 or n " 1{, the construction of GpGq uses more bits to express its rewards than what is used for the normal-form game G.

It is possible to avoid this phenomenon by instead using the state transitions of the Markov game to encode the action profile taken at each step, as was done in the proof of Theorem 3.2. The idea, which we sketch here, is to replace the game GpGq of Definition E.2 with the following game  $G^{1}pGq$ :

Definition E.9 (Alternative construction to Definition E.2). Given an m-player,  $n_0$ -action normal-form game G, we define the game  $G^1pGqpS;H;pA_iq_{iPr2s};P;pR_iq_{iPr2s};q$  as follows.

- The horizon of G is H n<sub>0</sub>.
- Let An<sub>0</sub>. The action spaces of agents 1;2;:::;m are given by A<sub>1</sub> A<sub>m</sub> rAs. The action space of agent m 1 is A<sub>m 1</sub>

so that  $|A_{m}|_{1} |Am/n$ .

We write A  $= {m \atop j1} A_j$  to denote the joint action space of the first m agents, and A  $:= {m \atop j1} A_j$  to denote the joint action space of all agents. Then  $|\overline{A}| A^m pmAq mA^{m-1}/n$ .

- The state space S is defined as follows. There are  $|\overline{A}|$  states, one for each action tuple a P  $\overline{A}$ . For each a P  $\overline{A}$ , we denote the corresponding state by  $s_a$ .
- For all h P r Hs, the reward to agent j P r m 1s given action profile a pa<sub>1</sub>;:::; $a_m$  1q at any state s P S is as follows: writing  $a_m$  1 pj<sup>1</sup>; $a_1^1$ q,

• At each step h PrHs, if action profile a  $\overline{PA}$  is taken, the game transitions to the state  $s_a$ .

Note that the number of states of  $G^1pGq$  is equal to  $|A| mn^m$  and so  $|G^1pGq| mn^m$ . As a result, if we were to use the game  $G^1pGq$  in place of GpGq in the proof of Theorem E.1, we would need to define  $n_0: tn^{1(pm)} q mu$  to ensure that  $|G^1pGq|/n$ , and so the condition T expp $^2tn\{mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(pm)} q mu\{m^2q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(pm)} q mu\{m^2q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(pm)} q mu\{m^2q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(qm)} q mu\{m^2q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expp $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^{1(qm)} q mu\{m^2q would be replaced by <math>T$  expe $^2tn^2q mu\{m^2q would be replaced by <math>T$  expe $^2tn^$ 

We expect that the construction of Definition E.2 can nevertheless still be modified to use Oplog1{q bits to express each reward in the Markov game G. In particular, one could introduce stochastic transitions to encode in the state of the Markov game a small number of random bits of the full action profile played at each step. We leave such an approach for future work.

## F. Equivalence between $\frac{gen;rnd}{i}$ and $p^{gen;det}q$

In this section we consider an alternate definition of the space  $\frac{\text{gen;rnd}}{i}$  of randomized general policies of player i, and show that it is equivalent to the one we gave in Section 2.

In particular, suppose we were to define a randomized general policy of agent i as a distribution over deterministic general policies of agent i: we write  $^{gen;rnd}$ :  $p^{gen;det}q_i$  to denote the space of such distributions. Moreover, write  $^{gen;rnd}$ :  $^{gen;rnd}$   $^{gen;rnd}$   $^{gen;rnd}$   $^{gen;rnd}$   $^{gen;rnd}$  to denote the space of product distributions over agents' deterministic policies. Our goal in this section is to show that policies in  $^{gen;rnd}$  are equivalent to those in  $^{gen;rnd}$  in the following sense: there is an embedding map  $^{gen;rnd}$ , not depending on the Markov game, so that the distribution of a trajectory drawn from any  $^{gen;rnd}$ , for any Markov game, is the same as the distribution of a trajectory drawn from Embpq (Fact F.2). Furthermore, Emb is surjective in the following sense: any policy  $^{gen;rnd}$  produces trajectories that are distributed identically to those of Embpq (and thus of ), for some  $^{gen;rnd}$  (Fact F.3). In Definition F.1 below, we define Emb.

Definition F.1. For j P rms and  $_j$  P  $_j^{\text{gen;rnd}}$ , define  $\text{Emb}_j p_j q P^{\text{gen;rnd}} p^{\text{gen;det}} q$  to put the following amount of mass on each  $_j$  P  $_j^{\text{gen;det}}$ :

$$pEmb_{j}p_{j}qqp_{j}q: " "^{H} _{j}p_{j;h}p_{j;h}; S_{h}q \mid_{j;h}; S_{h}q \mid_$$

Furthermore, for p<sub>1</sub>;:::;<sub>m</sub>q P<sup>gen;rnd</sup>, define Embpq pEmbp<sub>1</sub>q;:::;Embp<sub>m</sub>qq.

Note that, in the special case that  $_j$   $P^{\text{gen};d}_{j}$   $_e$ ,  $Emb_jp_jq$  is the point mass on  $_j$ .

Fact F.2 (Embedding equivalence). Fix a m-player Markov game G and, arbitrary policies  $_j$   $P^{\text{gen;rnd}}_{\phantom{gen;rnd}}$ . Then a trajectory drawn from the product policy  $p_1; ...;_{mq} P^{\text{gen;rnd}}_{\phantom{gen;rnd}}$  is distributed identically to a trajectory drawn from Embpq $P^{\text{gen;rnd}}$ .

The proof of Fact F.2 is provided in Section F.1. Next, we show that the mapping Emb is surjective in the following sense: Fact F.3 (Right inverse of  $Emb_j$ ). There is a mapping Fac: gen;rnd  $\tilde{N}$  gen;rnd so that for any Markov game G and any  $r P^{gen;rnd}$ , the distribution of a trajectory drawn from Facprq.

We will write  $Facppr_1; ...; r_mqq: pFac_1pr_1q; ...; Fac_mpr_mqq$ . Fact F.3 states that the policy Facprq maps, under Emb, to a policy in gen<sub>f</sub>rnd which is equivalent to r (in the sense that their trajectories are identically distributed for any Markov game).

An important consequence of Fact F.2 is that the expected reward (i.e., value) under any P<sup>gen;rnd</sup> is the same as that of Embpq. Thus given a Markov game, the induced normal-form game in which the players' pure action sets are <sup>gen;rnd</sup><sub>1</sub>; ...; <sup>gen;rnd</sup><sub>m</sub> is equivalent to the normal-form game in which the players' pure action sets are <sup>gen;det</sup><sub>1</sub>; ...; <sup>gen;det</sup><sub>m</sub>, in the following sense: for

any mixed strategy in the former, namely a product distributional policy  $P P p^{gen;rnd} q p^{gen;rnd} q$ , the policy  $E_P r E^{gen;rnd} q p^{gen;rnd} q$  is a mixed strategy in the latter which gives each player the same value as under P. (Note

that  $E_P$  rEmbpqs is indeed a product distribution since P is a product distribution and Emb factors into individual coordinates.) Furthermore, by Fact F.3, any distributional policy in  $g^{en;rnd}$  arises in this manner, for some  $P_1 P_2 P_3^{een;rnd} = P_3 P_3^{een;rnd} = P_3 P_3^{een;rnd} = P_$ 

#### F.1. Proofs of the equivalence

Proof of Fact F.2. Consider any trajectory  $ps_1; a_1; r_1; ...; s_H; a_H; r_H q$  consisting of a sequence of H states and actions and rewards for each of the m agents. Assume that  $r_{i;h} R_{i;h} ps; a_h q$  for all i;h (as otherwise has probability 0 under any policy). Write:

Then the probability of observing under is

where, per usual, j;h1 ps<sub>1</sub>; $a_{j;1};r_{j;1};...;s_{h1};a_{j;h1};r_{j;h1}q$ . Write p<sub>1</sub>;...;<sub>mq</sub> Embpq. The probability of observing under is

It is now straightforward to see from the definition of  $p_i$  in (24) that the quantities in (25) and (26) are equal.

Proof of Fact F.3. Fix a policy  $r_j$   $P_j^{\text{gen;rnd}}$   $p_j^{\text{gen;det}}q_{.j}$  We define  $\text{Fac}_j p_j q$  to be the policy  $_j$   $P_j^{\text{gen;rnd}}$ , which is defined as follows: for  $_{j;h1}$   $p_{s_{j;1};a_{j;1};r_{j;1};...;s_{j;h1};a_{j;h1};r_{j;h1}q$   $P_{j;h1}$ ,  $s_h$   $P_s$ , we have, for  $a_{j;h}$   $P_s$ ,

$$p_{j;h1} ; s ; qpa_{j;h}q = \frac{r_j^t P_j^{gen;det}}{r_j t_j P_j} : p_j; s_{j} ; qa_{j;g} @g/hu \\ = \frac{r_j^t P_j^{gen;det}}{r_j t_j P_j} : p_j; s_{j} ; s_{g} q_{a_{j};g} @g/h1u$$

If the denominator of the above expression is 0, then  $_{j}p_{j;h1}$ ;  $s_{h}q$  is defined to be an arbitrary distribution on  $pA_{j}q$ . (For concreteness, let us say that it puts all its mass on a fixed action in  $A_{j}$ .) Furthermore, for r P  $^{gen;rnd}r$ , define  $Facpq:pFac_{1}p_{1}q; r; Fac_{m}p_{m}qq P_{j}^{gen;rnd}$ .

Next, fix any r  $pr_1; ...; r_m q P^{gen} i_1^{rnd} gen; rnd$ . Let  $r_{m}^{Fac}$  Equation of trajectories under is the same as the distribution of trajectories drawn from .

So consider any trajectory  $ps_1; a_1; r_1; ...; s_H; a_H; r_H q$  consisting of a sequence of H states and actions and rewards for each of the m agents. Assume that  $r_{i;h}$   $R_{i;h}$   $ps_i; a_h q$  for all i;h (as otherwise has probability 0 under any policy). Write:

p: 
$$P_h ps_h _1 | s_h; a_h q: h1$$

Then the probability of observing under is

$$p^{"",H m} \\ p^{"",j;h}pa_{j;h}|_{j;h1};s_hq_{h1j1} \\ p^{"m",H} \\ r_j t_j P^{gen;d_t}_{j_e}:_{j}p_{j;g};s_gqa_{j;g} @g/hu \\ p^{gen;d_t} t_j^{r}P^{gen;d_t}_{e}:_{j}p_{j;g};s_gqa_{j;g} @g/h1up^{"}_{j}t_j \\ P^{gen;d_t}_{e} r_{j}p_{j;g};s_gqa_{j;g} @g/Hu ; \\ p^{gen;d_t}_{j} r_{j}p_{j;g};s_gqa_{j;g} @g/Hu ; \\ p^{gen;d_t}_{j} r_{j}p_{j;g};s_gqa_{j;g} gg/Hu ; \\ p^{g$$

which is equal to the probability of observing under r.