
When is Mean-Field Reinforcement Learning Tractable and Relevant?

Anonymous Authors¹

Abstract

Mean-field reinforcement learning has become a popular theoretical framework for efficiently approximating large-scale multi-agent reinforcement learning (MARL) problems exhibiting symmetry. However, questions remain regarding the applicability of mean-field approximations: in particular, their approximation accuracy of real-world systems and conditions under which they become computationally tractable. We establish explicit finite-agent bounds for how well the MFG solution approximates the true N -player game for two popular mean-field solution concepts. Furthermore, for the first time, we establish explicit lower bounds indicating that MFGs are poor or uninformative at approximating N -player games assuming only Lipschitz dynamics and rewards. Finally, we analyze the computational complexity of solving MFGs with only Lipschitz properties and prove that they are in the class of PPAD-complete problems conjectured to be intractable, similar to general sum N player games. Our theoretical results underscore the limitations of MFGs and complement and justify existing work by proving difficulty in the absence of common theoretical assumptions.

1. Introduction

Multi-agent reinforcement learning (MARL) finds numerous impactful applications in the real world (Shavandi & Khedmati, 2022; Wiering, 2000; Samvelyan et al., 2019; Rashedi et al., 2016; Matignon et al., 2007; Mao et al., 2022). Despite the urgent need in practice, MARL remains a fundamental challenge, especially in the setting with large numbers of agents due to the so-called “curse of many agents” (Wang et al., 2020).

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the Workshop on Foundations of Reinforcement Learning and Control at the International Conference on Machine Learning (ICML). Do not distribute.

Mean-field games (MFG), a theoretical framework first proposed by Lasry & Lions (2007) and Huang et al. (2006), permits the theoretical study of such large-scale games by introducing mean-field simplification. Under certain assumptions, the mean-field approximation leads to efficient algorithms for the analysis of a particular type of N -agent competitive game where there are symmetries between players and when N is large. Such games appear widely in for instance auctions (Iyer et al., 2014), and cloud resource management (Mao et al., 2022). For the mean-field analysis, the game dynamics with N -players must be *symmetric* (i.e., each player must be exposed to the same rules) and *anonymous* (i.e., the effect of each player on the others should be permutation invariant). Under this simplification, works such as (Perrin et al., 2020; Anahtarci et al., 2022; Guo et al., 2019; Pérolat et al., 2022; Xie et al., 2021) and many others have analyzed reinforcement learning (RL) algorithms in the MFG limit $N \rightarrow \infty$ to obtain a tractable approximation of many agent games, providing learning guarantees under various structural assumptions.

Being a simplification, MFG formulations should ideally satisfy two desiderata: (1) they should be *relevant*, i.e., they are good approximations of the original MARL problem and (2) they should be *tractable*, i.e., they are at least easier than solving the original MARL problem. In this work, we would like to understand the extent to which MFGs satisfy these two requirements, and we aim to answer two natural questions that remain understudied:

- *When are MFGs good approximations of the finite player games, when are they not?* In particular, are polynomially many agents always sufficient for mean-field approximation to be effective?
- *Is solving MFGs always computationally tractable, or more tractable than directly solving the N -player game?* In particular, can MFGs be solved in polynomial or pseudo-polynomial time?

1.1. Related Work

Mean-field RL has been studied in various mathematical settings. In this work, we focus on two popular formulations in particular: stationary mean-field games (Stat-MFG, see e.g. (Anahtarci et al., 2022; Guo et al., 2019)) and finite-horizon

MFG (FH-MFG, see e.g. (Perrin et al., 2020; Pérolat et al., 2022)). In the Stat-MFG setting the objective is to find a stationary policy that is optimal with respect to its induced stationary distribution, while in the FH-MFG setting, a finite-horizon reward is considered with a time-varying policy and population distribution.

Existing results on MFG relevance/approximation. The approximation properties of MFGs have been explored by several works in literature, as summarized in Table 1. Finite-agent approximation bounds have been widely analyzed in the case of stochastic mean-field differential games (Carmona & Delarue, 2013; Carmona et al., 2018), albeit in the differential setting and without explicit lower bounds. Recent works (Anahtarci et al., 2022; Cui & Koepl, 2021) have established that Stat-MFG Nash equilibria (Stat-MFG-NE) asymptotically approximate the NE of N -player symmetric dynamic games under continuity assumptions. The result by Saldi et al. (2018), as the basis of subsequent proofs, shows asymptotic convergence for a large class of MFG variants and only requires continuity of dynamics and rewards as well as minor technical assumptions such as compactness and a form of local Lipschitz continuity. However, such asymptotic convergence guarantees leave the question unanswered if the MFG models are realistic in real-world games. Many games such as traffic systems, financial markets, etc. naturally exhibit large N , however, if N must be astronomically large for good approximation, the real-world impact of the mean-field analysis will be limited. Recently, (Yardim et al., 2023b) provided finite-agent approximation bounds of a special class of stateless MFG, which assumes no state dynamics. We complement existing work on approximation properties of both Stat-MFG and FH-MFG by providing explicit upper and lower bounds for approximation.

Existing results on MFG tractability. The tractability of solving MFGs as a proxy for MARL has been also heavily studied in the RL community under various classes of structural assumptions. Since finding approximate Nash equilibria for normal form games is PPAD-complete, a class believed to be computationally intractable (Daskalakis et al., 2009; Chen et al., 2009), solving the mean-field approximation in many cases can be a tractable alternative. We summarize recent work for computationally (or statistically) solving the two types of MFGs below, with an in-depth comparison also provided in Table 2.

For Stat-MFG, under a contraction assumption RL algorithms such as Q-learning (Zaman et al., 2023; Anahtarci et al., 2022), policy mirror ascent (Yardim et al., 2023a), policy gradient methods (Guo et al., 2022a), soft Q-learning (Cui & Koepl, 2021) and fictitious play (Xie et al., 2021) have been shown to solve Stat-MFG with statistical and computational efficiency. However, all of these guarantees

require the game to be heavily regularized as pointed out in (Cui & Koepl, 2021; Yardim et al., 2023a), inducing a non-vanishing bias on the computed Nash. Moreover, in some works the population evolution is also implicitly required to be contractive under all policies (see e.g. (Guo et al., 2019; Yardim et al., 2023a)), further restricting the analysis to sufficiently smooth games. While (Guo et al., 2022b) has proposed a method that guarantees convergence to MFG-NE under differentiable dynamics, the algorithm converges only when initialized sufficiently close to the solution. To the best of our knowledge, there are neither RL algorithms that work without regularization nor evidence of difficulty in the absence of such strong assumptions: we complement the line of work by showing that unless dynamics are sufficiently smooth, Stat-MFG is both computationally intractable and a poor approximation.

A separate line of work analyzes the finite horizon problem. In this case, when the dynamics are population-independent and the payoffs are monotone the problem is known to be tractable. Algorithms such as fictitious play (Perrin et al., 2020) and mirror descent (Pérolat et al., 2022) have been shown to converge to Nash in corresponding continuous-time equations. Recent work has also focused on the statistical complexity of the finite-horizon problem in very general FH-MFG problems (Huang et al., 2023), however, the algorithm proposed is in general computationally intractable. In terms of computational tractability and the approximation properties, our work complements these results by demonstrating that (1) when dynamics depend on the population as well an exponential approximation lower bound exists, and (2) in the absence of monotonicity, the FH-MFG is provably as difficult as solving an N -player game.

1.2. Our Contribution

In this work, we formalize and provide answers to the two aforementioned fundamental questions, first focusing on the approximation properties of MFG in Section 3 and later on the computational tractability of MFG in Section 4. Our contributions are summarized as follows.

Firstly, we introduce explicit finite-agent approximation bounds for finite horizon and stationary MFGs (Table 1) in terms of exploitability in the finite agent game. In both cases, we prove explicit upper bounds which quantify how many agents a symmetric game must have to be well-approximated by the MFG, which has been absent in the literature to the best of our knowledge. Our approximation results only require a minimal Lipschitz continuity assumption of the transition kernel and rewards. For FH-MFG, we prove a $\mathcal{O}\left(\frac{(1-L^H)H^2}{(1-L)\sqrt{N}}\right)$ upper bound for the exploitability where L is the Lipschitz modulus of the population evolution operator: the upper bound exhibits an exponential dependence on the horizon H . For the Stat-MFG we

show that a $\mathcal{O}\left(\frac{(1-\gamma)^{-3}}{\sqrt{N}}\right)$ approximation bound can be established, but only if the population evolution dynamics are non-expansive. Next, for the first time, we establish explicit lower bounds for the approximation proving the shortcomings of the upper bounds are fundamental. For the FH-MFG, we show that unless $N \geq \Omega(2^H)$, an exploitability linear in horizon H is unavoidable when deploying the MFG solution to the N player game: hence in general the MFG equilibrium becomes irrelevant quickly as the problem horizon increases. For Stat-MFG we establish an $\Omega(N^{\log_2 \gamma})$ lower bound when the population dynamics are not restricted to non-expansive population operators, showing that a large discount factor γ also rapidly deteriorates the approximation efficiency. Our lower bounds indicate that in the worst case, the number of agents required for the approximation can grow exponentially in the problem parameters, demonstrating the limitations of the MFG approximation.

Finally, from the computational perspective, we establish that both finite-horizon and stationary MFGs can be PPAD-complete problems in general, even when restricted to certain simple subclasses (Table 2). This shows that both MFG problems are in general as hard as finding a Nash equilibrium of N -player general sum games. Furthermore, our results imply that unless PPAD=P there are no polynomial time algorithms for solving FH-MFG and Stat-MFG, a result indicating computational intractability.

2. Mean-Field Games: Definitions, Solution Concepts

Notation. Throughout this work, we assume \mathcal{S}, \mathcal{A} are finite sets. For a finite set \mathcal{X} , $\Delta_{\mathcal{X}}$ denotes the set of probability distributions on \mathcal{X} . The norm used will not fundamentally matter for our results, we choose to equip $\Delta_{\mathcal{S}}, \Delta_{\mathcal{A}}$ with the norm $\|\cdot\|_1$. We define the set of Markov policies $\Pi := \{\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}\}$, $\Pi_H := \{\{\pi_h\}_{h=0}^{H-1} : \pi_h \in \Pi, \forall h\}$ and $\Pi_H^N := \{\{\pi_h^i\}_{h=0, i=0}^{H-1, N} : \pi_h^i \in \Pi, \forall h\}$. For policies $\pi, \pi' \in \Pi$ denote $\|\pi - \pi'\|_1 = \sup_{s \in \mathcal{S}} \|\pi(\cdot|s) - \pi'(\cdot|s)\|_1$. We denote $d(x, y) := \mathbb{1}_{\{x \neq y\}}$ for x, y in \mathcal{A} or \mathcal{S} . For $\pi \in \Pi^N, \pi' \in \Pi$, we define $(\pi', \pi^{-i}) \in \Pi^N$ as the policy profile where the i -th policy has been replaced by π' . Likewise, for $\pi \in \Pi_H^N, \pi' \in \Pi_H$, we denote by $(\pi', \pi^{-i}) \in \Pi_H^N$ the policy profile where the i -th player's policy has been replaced by π' . For any $N \in \mathbb{N}_{\geq 0}$, $[N] := \{1, \dots, N\}$.

MFGs introduce a dependence on the population distribution over states of the rewards and dynamics. We will strictly consider Lipschitz continuous rewards and dynamics, which is a common assumption in literature (Guo et al., 2019; Anaharci et al., 2022; Yardim et al., 2023a; Xie et al., 2021), formalized below.

Definition 2.1 (Lipschitz dynamics, rewards). For some $L \geq 0$, we define the set of L -Lipschitz reward functions

and state transition dynamics as

$$\begin{aligned} \mathcal{R}_L &:= \left\{ R : \mathcal{S} \times \mathcal{A} \times \Delta_{\mathcal{S}} \rightarrow [0, 1] : \right. \\ &\quad \left. |R(s, a, \mu) - R(s, a, \mu')| \leq L \|\mu - \mu'\|_1, \forall s, a, \mu, \mu' \right\}, \\ \mathcal{P}_L &:= \left\{ P : \mathcal{S} \times \mathcal{A} \times \Delta_{\mathcal{S}} \rightarrow \Delta_{\mathcal{S}} : \right. \\ &\quad \left. \|P(s, a, \mu) - P(s, a, \mu')\|_1 \leq L \|\mu - \mu'\|_1, \forall s, a, \mu, \mu' \right\}. \end{aligned}$$

Moreover, we define the set of Lipschitz rewards and dynamics as $\mathcal{R} := \bigcup_{L \geq 0} \mathcal{R}_L$, $\mathcal{P} := \bigcup_{L \geq 0} \mathcal{P}_L$ respectively.

We note that there are interesting MFGs with non-Lipschitz dynamics and rewards, however, even the existence of Nash is not guaranteed in this case. Lipschitz continuity is a minimal assumption under which solutions to MFG always exist, and as our aim is to prove lower bounds and difficulty we will adopt this assumption. Solving MFG with non-Lipschitz dynamics is more challenging than Lipschitz continuous MFG (the latter being a subset of the former), hence our difficulty results will apply.

Operators. We will define the useful population operators $\Gamma_P : \Delta_{\mathcal{S}} \times \Pi \rightarrow \Delta_{\mathcal{S}}$, $\Gamma_P^H : \Delta_{\mathcal{S}} \times \Pi \rightarrow \Delta_{\mathcal{S}}$, and $\Lambda_P^H : \Delta_{\mathcal{S}} \times \Pi_H \rightarrow \Delta_{\mathcal{S}}$ as

$$\begin{aligned} \Gamma_P(\mu, \pi) &:= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \mu(s) \pi(a|s) P(\cdot|s, a, \mu), \\ \Gamma_P^H(\mu, \pi) &:= \underbrace{\Gamma_P(\dots \Gamma_P(\Gamma_P(\mu, \pi), \pi) \dots)}_{H \text{ times}}, \\ \Lambda_P^H(\mu_0, \boldsymbol{\pi}) &:= \left\{ \underbrace{\Gamma_P(\dots \Gamma_P(\Gamma_P(\mu_0, \pi_0), \pi_1) \dots, \pi_{h-1})}_{h \text{ times}} \right\}_{h=0}^{H-1} \end{aligned}$$

for all $n \in \mathbb{N}_{>0}, \pi \in \Pi, \boldsymbol{\pi} = \{\pi_h\}_{h=0}^{H-1} \in \Pi_H, P \in \mathcal{P}, \mu_0 \in \Delta_{\mathcal{S}}$.

Finally, we will need the following Lipschitz continuity result for the Γ_P operator.

Lemma 2.2 (Lemma 3.2 of (Yardim et al., 2023a)). *Let $P \in \mathcal{P}_{K_\mu}$ for $K_\mu > 0$ and*

$$\begin{aligned} K_s &:= \sup_{\substack{s, s' \\ a, \mu}} \|P(s, a, \mu) - P(s', a, \mu)\|_1, \\ K_a &:= \sup_{\substack{a, a' \\ s, \mu}} \|P(s, a, \mu) - P(s, a', \mu)\|_1. \end{aligned}$$

Then it holds for all $\mu, \mu' \in \Delta_{\mathcal{S}}, \pi, \pi' \in \Pi$ that:

$$\begin{aligned} \|\Gamma_P(\mu, \pi) - \Gamma_P(\mu', \pi')\|_1 &\leq L_{pop, \mu} \|\mu - \mu'\|_1 \\ &\quad + \frac{K_a}{2} \|\pi - \pi'\|_1, \end{aligned}$$

$\forall \pi, \pi' \in \Pi, \mu, \mu' \in \Delta_{\mathcal{S}}$, and $L_{pop, \mu} := (K_\mu + \frac{K_s}{2} + \frac{K_a}{2})$.

Tractability and Relevance of MF-RL

Work	MFG type	Key Assumptions	Approximation Rate (in Exploitability)
Carmona et al., 2013	Other ^a	Affine drift, Lip. derivatives	$\mathcal{O}(N^{-1/(d+4)})$ (d : dim. of state space)
Saldi et al., 2018	Other ^b	Continuity	$o(1)$ (convergence as $N \rightarrow \infty$)
Anahtarci et al., 2022	Stat-MFG	Lip. P, R + Reg. + Contractive Γ_P	$o(1)$ (convergence as $N \rightarrow \infty$)
Cui & Koepl, 2021	Stat-MFG	Continuity	$o(1)$ (convergence as $N \rightarrow \infty$)
Yardim et al., 2023b	Other ^c	Lip. P, R	$\mathcal{O}(1/\sqrt{N})$
Theorem 3.2	FH-MFG	Lip. P, R	$\mathcal{O}\left(\frac{H^2(1-L^H)}{(1-L)\sqrt{N}}\right)$, L Lip. modulus of Γ_P
Theorem 3.3	FH-MFG	Lip. P, R	$\Omega(H)$ unless $N \geq \Omega(2^H)$
Theorem 3.5	Stat-MFG	Lip. P, R + Non-expansive Γ_P	$\mathcal{O}((1-\gamma)^{-3}/\sqrt{N})$
Theorem 3.6	Stat-MFG	Lip. P, R	$\Omega(N^{-\log_2 \gamma^{-1}})$

Table 1. Selected approximation results for MFG. Notes: ^a stochastic differential MFG, ^b infinite-horizon discounted setting with non-stationary policies, ^c stateless/static MFG setting. Lip.=Lipschitz, Reg.=non-vanishing regularization required.

Work	MFG Type	Key Assumptions	Iteration/Sample Complexity result
Anahtarci et al., 2022	Stat-MFG	Lip. P, R + Reg. + Contractive Γ_P	$\tilde{\mathcal{O}}(\varepsilon^{-4 A })$ samples, $\mathcal{O}(\log \varepsilon^{-1})$ iterations
Geist et al., 2022	Other ^a	Concave potential	$\mathcal{O}(\varepsilon^{-2})$ iterations
Perrin et al., 2020	FH-MFG	Monotone R, μ -independent P	$\mathcal{O}(\varepsilon^{-1})$ (continuous time analysis)
Pérolat et al., 2022	FH-MFG	Monotone R, μ -independent P	$\mathcal{O}(\varepsilon^{-1})$ (continuous time analysis)
Zaman et al., 2023	Stat-MFG	Lip. P, R + Reg. + Contractive Γ_P	$\mathcal{O}(\varepsilon^{-4})$ samples
Cui & Koepl, 2021	Stat-MFG	Lip. P, R + Reg.	$\mathcal{O}(\log \varepsilon^{-1})$ iterations
Yardim et al., 2023b	Other ^b	Monotone and Lip. R	$\tilde{\mathcal{O}}(\varepsilon^{-2})$ samples (N -player)
Yardim et al., 2023a	Stat-MFG	Lip. P, R + Reg. + Contractive Γ_P	$\tilde{\mathcal{O}}(\varepsilon^{-2})$ samples (N -player)
Theorem 4.9	Stat-MFG	Lip. P, R	PPAD-complete
Theorem 4.12	FH-MFG	Lip. P, R + μ -independent P	PPAD-complete
Theorem 4.14	FH-MFG	Linear P, R + μ -independent P	PPAD-complete

Table 2. Selected results for computing MFG-NE from literature. In the assumptions column, contractive Γ_P indicates that for all $\pi \in \Pi$, $\Gamma_P(\cdot, \pi)$ is a contraction, and regularization indicates that a non-vanishing bias is present. Notes: ^a infinite-horizon, population dependence through the discounted state distribution. ^b stateless/static MFG. Lip.=Lipschitz, Reg.=non-vanishing regularization required.

In particular, in our settings, Lemma 2.2 indicates that Γ_P is always Lipschitz continuous if $P \in \mathcal{P}$, a property which will become significant for approximation analysis.

We will be interested in two classes of MFG solution concepts that lead to different analyses: infinite horizon stationary MFG Nash equilibrium (Stat-MFG-NE) and finite horizon MFG Nash equilibrium (FH-MFG-NE). The first problem widely studied in literature is the stationary MFG equilibrium problem, see for instance (Anahtarci et al., 2022; Yardim et al., 2023a; Guo et al., 2019; 2022a; Xie et al., 2021). We formalize this solution concept below.

Definition 2.3 (Stat-MFG). A stationary MFG (Stat-MFG) is defined by the tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ for Lipschitz dynamics and rewards $P \in \mathcal{P}$, $R \in \mathcal{R}$, discount factor $\gamma \in (0, 1)$. For any $(\mu, \pi) \in \Delta_{\mathcal{S}} \times \Pi$, we define the γ -discounted

infinite horizon expected reward as

$$V_{P,R}^{\gamma}(\mu, \pi) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, \mu) \middle| \begin{matrix} s_0 \sim \mu, & a_t \sim \pi(s_t) \\ s_{t+1} \sim P(s_t, a_t, \mu) \end{matrix} \right].$$

A policy-population pair $(\mu^*, \pi^*) \in \Delta_{\mathcal{S}} \times \Pi$ is called a Stat-MFG Nash equilibrium if the two conditions hold:

$$\begin{aligned} \text{Stability:} & \quad \mu^* = \Gamma_P(\mu^*, \pi^*), \\ \text{Optimality:} & \quad V_{P,R}^{\gamma}(\mu^*, \pi^*) = \max_{\pi \in \Pi} V_{P,R}^{\gamma}(\mu^*, \pi). \end{aligned} \tag{Stat-MFG-NE}$$

The second MFG concept that we will consider has a finite time horizon, and is also common in literature (Perolat et al., 2015; Perrin et al., 2020; Laurière et al., 2022; Huang et al., 2023). In this case, the population distribution is permitted to vary over time, and the objective is to find an optimal non-stationary policy with respect to the population distribution it induces. We formalize this problem and the corresponding solution concept below.

Definition 2.4 (FH-MFG). A finite horizon MFG problem (FH-MFG) is determined by the tuple $(\mathcal{S}, \mathcal{A}, H, P, R, \mu_0)$ where $H \in \mathbb{Z}_{>0}$, $P \in \mathcal{P}$, $R \in \mathcal{R}$, $\mu_0 \in \Delta_{\mathcal{S}}$. For $\pi = \{\pi_h\}_{h=0}^H \in \Pi_H$, $\mu = \{\mu_h\}_{h=0}^{H-1} \in \Delta_{\mathcal{S}}^H$, define the expected reward and exploitability as

$$V_{P,R}^H(\mu, \pi) := \mathbb{E} \left[\sum_{h=0}^{H-1} R(s_h, a_h, \mu_h) \middle| \begin{array}{l} s_0 \sim \mu_0, a_h \sim \pi_h(s_h) \\ s_{h+1} \sim P(s_h, a_h, \mu_h) \end{array} \right],$$

$$\mathcal{E}_{P,R}^H(\pi) := \max_{\pi' \in \Pi_H} V_{P,R}^H(\Lambda_P^H(\mu_0, \pi), \pi') - V_{P,R}^H(\Lambda_P^H(\mu_0, \pi), \pi).$$

Then, the FH-MFG Nash equilibrium is defined as:

$$\text{Policy } \pi^* = \{\pi_h^*\}_{h=0}^{H-1} \in \Pi_H \text{ such that}$$

$$\mathcal{E}_{P,R}^H(\{\pi_h^*\}_{h=0}^{H-1}) = 0. \quad (\text{FH-MFG-NE})$$

3. Approximation Properties of MFG

As established in literature, the reason the FH-MFG and Stat-MFG problems are studied is the fact that they can approximate the NE of certain symmetric games with N players, establishing the main relevance of the formulations in the real world. Such results are summarized in Table 1.

In this section, we study how efficient this convergence is and also related lower bounds. For these purposes, we first define the corresponding *finite-player* game of each mean-field game problem: to avoid confusion, we call these games *symmetric anonymous dynamic games* (SAG). Afterwards, for each solution concept, we will first establish (1) an upper bound on the approximation error (i.e. the exploitability) due to the mean-field, and (2) a lower bound demonstrating the worst-case rate. We will present the main outlines of proofs, and postpone computation-intensive derivations to the supplementary material of the paper.

3.1. Approximation Analysis of FH-MFG

Firstly, we define the finite-player game that is approximately solved by the FH-MFG-NE.

Definition 3.1 (N -FH-SAG). An N -player finite horizon SAG (N -FH-SAG) is determined by the tuple $(N, \mathcal{S}, \mathcal{A}, H, P, R, \mu_0)$ such that $N \in \mathbb{Z}_{>0}$, $H \in \mathbb{Z}_{>0}$, $P \in \mathcal{P}$, $R \in \mathcal{R}$, $\mu_0 \in \Delta_{\mathcal{S}}$. For any $\pi = \{\pi_h^i\}_{h=0, \dots, H-1, i \in [N]} \in \Pi_H^N$, we define the expected mean reward and exploitability of player i as

$$J_{P,R}^{H,N,(i)}(\pi) := \mathbb{E} \left[\sum_{h=0}^{H-1} R(s_h^i, a_h^i, \hat{\mu}_h) \middle| \begin{array}{l} \forall j: s_0^j \sim \mu_0, a_h^j \sim \pi_h^j(s_h^j) \\ \hat{\mu}_h := \frac{1}{N} \sum_j e^{s_h^j} \\ s_{h+1}^j \sim P(s_h^j, a_h^j, \hat{\mu}_h) \end{array} \right],$$

$$\mathcal{E}_{P,R}^{H,N,(i)}(\pi) := \max_{\pi' \in \Pi_H^N} J_{P,R}^{H,N,(i)}(\pi', \pi^{-i}) - J_{P,R}^{H,N,(i)}(\pi).$$

Then, the N -FH-SAG Nash equilibrium is defined as:

$$N\text{-tuple of policies } \{\pi_h^{(i),*}\}_{h=0}^{H-1} \in \Pi_H^N \text{ such that}$$

$$\forall i: \mathcal{E}_{P,R}^{H,N,(i)}(\{\pi_h^*\}_{h=0}^{H-1}) = 0. \quad (N\text{-FH-SAG-NE})$$

If instead $\mathcal{E}_{P,R}^{H,N,(i)}(\pi) \leq \delta$ for all i , then π is called a δ - N -FH-SAG Nash equilibrium.

The above definition corresponds to a real-world problem as the function $J_{P,R}^{H,N,(i)}$ expresses the expected total payoff of each player: hence a δ - N -MFG-NE is a Nash equilibrium of a concrete N -player game in the traditional game theoretical sense. Also, note that now in the definition transition probabilities and rewards depend on $\hat{\mu}_h$ which is the $\mathcal{F}(\{s_h^i\}_i) = \mathcal{F}_h$ -measurable random vector of the empirical state distribution at time h of all agents.

Firstly, we provide a positive result well-known in literature: the N -FH-SAG is approximately solved by the FH-MFG-NE policy. Unlike some past works, we establish an explicit rate of convergence in terms of N and problem parameters.

Theorem 3.2 (Approximation of N -FH-SAG). *Let $(\mathcal{S}, \mathcal{A}, H, P, R, \mu_0)$ be a FH-MFG with $P \in \mathcal{P}$, $R \in \mathcal{R}$ and with a FH-MFG-NE $\pi^* \in \Pi_H$, and for any $N \in \mathbb{N}_{>0}$ let $\pi_N^* := \underbrace{(\pi^*, \dots, \pi^*)}_{N \text{ times}} \in \Pi_H^N$. Let $L > 0$ be the Lipschitz constant of Γ_P in μ , and let $\mathcal{G}_N := (N, \mathcal{S}, \mathcal{A}, H, P, R, \mu_0)$ be the corresponding N -player game. Then:*

1. *If $L = 1$, then for all $i \in [N]$, $\mathcal{E}_{P,R}^{H,N,(i)}(\pi_N^*) \leq \mathcal{O}(\frac{H^3}{\sqrt{N}})$, that is, π_N^* is a $\mathcal{O}(\frac{H^3}{\sqrt{N}})$ -NE of \mathcal{G}_N .*
2. *If $L \neq 1$, then for all $i \in [N]$, $\mathcal{E}_{P,R}^{H,N,(i)}(\pi_N^*) \leq \mathcal{O}\left(\frac{H^2(1-L^H)}{(1-L)\sqrt{N}}\right)$, that is, π_N^* is a $\mathcal{O}\left(\frac{H^2(1-L^H)}{(1-L)\sqrt{N}}\right)$ -NE of \mathcal{G}_N .*

Γ_P in Theorem 3.2 is always L -Lipschitz in μ for some L by Lemma 2.2. When $L > 1$, the upper bound $\mathcal{O}((1+L^H)H^2/\sqrt{N})$ has an exponential dependence on the Lipschitz constant of the operator Γ_P . However, for games with longer horizons, the upper bound might require an unrealistic amount of agents N to guarantee a good approximation due to the exponential dependency. Next, we establish a worst-case result demonstrating that this is not avoidable without additional assumptions.

Theorem 3.3 (Approximation lower bound for N -FH-SAG). *There exists \mathcal{S}, \mathcal{A} and $P \in \mathcal{P}_8, R \in \mathcal{R}_2, \mu_0 \in \Delta_{\mathcal{S}}$ such that the following hold:*

1. *For each $H > 0$, the FH-MFG defined by $(\mathcal{S}, \mathcal{A}, H, P, R, \mu_0)$ has a unique solution π_H^* (up to modifications on zero-probability sets),*
2. *For any $H, h > 0$, in the N -FH-SAG it holds that $\mathbb{E}_H[\|\hat{\mu}_h - \Lambda_P^H(\mu_0, \pi_H^*)_h\|_1] \geq \Omega\left(\min\left\{1, \frac{2^H}{\sqrt{N}}\right\}\right)$.*
3. *For any $H, N > 0$ either $N \geq \Omega(2^H)$, or for each player $i \in [N]$ it holds that $\mathcal{E}_{P,R}^{H,N,(i)}(\pi_H^*, \dots, \pi_H^*) \geq \Omega(H)$.*

This result shows that without further assumptions, the FH-MFG solution might suffer from exponential exploitability in H in the N -player game. In such cases, to avoid the concrete N -player game from deviating from the mean-field behavior too fast, either H must be small or P must be sufficiently smooth in μ . We note that the typical assumption in the finite-horizon setting that $P \in \mathcal{P}_0$ (see e.g. (Perrin et al., 2020; Geist et al., 2022)) avoids this lower bound since in this case $\Gamma_P(\cdot, \pi)$ is simply multiplication by a stochastic matrix which is always non-expansive ($L = 1$). We also note at the expense of simplicity a stronger counter-example inducing exploitability $\Omega(H)$ unless $N \geq \Omega((L - \epsilon)^H)$ for all $\epsilon > 0$ can be constructed, where $P \in \mathcal{P}_L$.

A remark. The proof of Theorem 3.3 in fact suggests that for finite N and large horizon H , there exists a time-homogenous policy $\bar{\pi}^* \in \Pi$ different than the FH-MFG solution such that for $\bar{\pi}_H^* := \{\bar{\pi}^*\}_{h=0}^{H-1} \in \Pi_H$, the time-averaged exploitability of $\bar{\pi}_H^*$ is small: $\forall i \in [N] : H^{-1} \mathcal{E}_{P,R}^{H,N,(i)}(\bar{\pi}_H^*, \dots, \bar{\pi}_H^*) \leq \mathcal{O}(H^{-1} \log_2 N)$.

3.2. Approximation Analysis of Stat-MFG

Similarly, we introduce the N -player game corresponding to the Stat-MFG solution concept.

Definition 3.4 (N -Stat-SAG). An N -player stationary SAG (N -Stat-SAG) problem is defined by the tuple $(N, \mathcal{S}, \mathcal{A}, P, R, \gamma)$ for Lipschitz dynamics and rewards $P \in \mathcal{P}, R \in \mathcal{R}$, discount factor $\gamma \in (0, 1)$. For any $(\mu, \pi) \in \Delta_{\mathcal{S}} \times \Pi^N$, the N -player γ -discounted infinite horizon expected reward is defined as:

$$J_{P,R}^{\gamma,N,(i)}(\mu, \pi) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^i, a_t^i, \hat{\mu}_t) \middle| \begin{array}{l} a_t^i \sim \pi^j(s_t^i), \hat{\mu}_t := \frac{\sum_j e^{s_t^j}}{N} \\ s_0^i \sim \mu, s_{t+1}^i \sim P(s_t^i, a_t^i, \hat{\mu}_t) \end{array} \right].$$

A policy profile-population pair $(\mu^*, \pi^*) \in \Delta_{\mathcal{S}} \times \Pi^N$ is called an N -Stat-SAG Nash equilibrium if:

$$J_{P,R}^{\gamma,N,(i)}(\mu^*, \pi^*) = \max_{\pi \in \Pi} J_{P,R}^{\gamma,N,(i)}(\mu^*, (\pi, \pi^{*-i})). \quad (N\text{-Stat-SAG-NE})$$

If $J_{P,R}^{\gamma,N,(i)}(\mu^*, \pi^*) \geq \max_{\pi \in \Pi} J_{P,R}^{\gamma,N,(i)}(\mu^*, (\pi, \pi^{*-i})) - \delta$, then we call μ^*, π^* a δ - N -Stat-SAG Nash equilibrium.

Theorem 3.5 (Approximation of N -Stat-SAG). Let $(\mathcal{S}, \mathcal{A}, H, P, R, \gamma)$ be a Stat-MFG and $(\mu^*, \pi^*) \in \Delta_{\mathcal{S}} \times \Pi$ be a corresponding Stat-MFG-NE. Furthermore, assume that $\Gamma_P(\cdot, \pi)$ is non-expansive in the ℓ_1 norm for any π , that is, $\|\Gamma_P(\mu, \pi) - \Gamma_P(\mu', \pi)\|_1 \leq \|\mu - \mu'\|_1$. Then, $(\mu^*, \pi^*) \in \Delta_{\mathcal{S}} \times \Pi^N$ is a $\mathcal{O}\left(\frac{1}{\sqrt{N}}\right)$ Nash equilibrium for the N -player game where $\pi_N^* := (\pi^*, \dots, \pi^*)$, that is, for all i ,

$$J_{P,R}^{\gamma,N,(i)}(\mu^*, \pi_N^*) \geq \max_{\pi \in \Pi} J_{P,R}^{\gamma,N,(i)}(\mu^*, (\pi, \pi_N^{*-i})) - \mathcal{O}\left(\frac{(1-\gamma)^{-3}}{\sqrt{N}}\right).$$

We also establish an approximation lower bound for the N -Stat-SAG. In this case, the question is if the non-expansive Γ_P assumption is necessary for the optimal $\mathcal{O}(1/\sqrt{N})$ rate. The below results affirm this: in for Stat-MFG-NE with expansive Γ_P , we suffer from an exploitability of $\omega(1/\sqrt{N})$ in the N -agent case.

Theorem 3.6 (Lower bound for N -Stat-SAG). For any $N \in \mathbb{N}_{>0}, \gamma \in (1/\sqrt{2}, 1)$ there exists \mathcal{S}, \mathcal{A} with $|\mathcal{S}| = 6, |\mathcal{A}| = 2$ and $P \in \mathcal{P}_7, R \in \mathcal{R}_3$ such that:

1. The Stat-MFG $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ has a unique NE μ^*, π^* ,
2. For any N and $\pi_N^* := (\pi^*, \dots, \pi^*) \in \Pi^N$, it holds that $J_{P,R}^{\gamma,N,(i)}(\pi_N^*) \leq \max_{\pi} J_{P,R}^{\gamma,N,(i)}(\pi, \pi_N^{*-i}) - \Omega(N^{-\log_2 \gamma^{-1}})$.

The result above shows that unless the relevant Γ_P operator is contracting in some potential, in general, the exploitability of the Stat-MFG-NE in the N -player game might be very large unless the effective horizon $(1-\gamma)^{-1}$ is small. Hence, in these cases, the mean-field Nash equilibrium might be uninformative regarding the true NE of the N player game. In the case of Stat-MFG, our lower bound is even stronger in the sense that the exploitability no longer decreases with $\mathcal{O}(1/\sqrt{N})$ for large γ . For a sufficiently long effective horizon $(1-\gamma)^{-1}$ and large enough Lipschitz constant L , the rate in terms of N can be arbitrarily slow. Furthermore, if we take the ergodic limit $\gamma \rightarrow 1$, we will observe a non-vanishing exploitability $\Omega(1)$ for all finite N .

4. Computational Tractability of MFG

The next fundamental question for mean-field reinforcement learning will be whether it is always computationally easier than finding an equilibrium of a N -player general sum normal form game. We focus on the computational aspect of solving mean-field games in this section, and not statistical uncertainty: we assume we have full knowledge of the MFG dynamics. We will show that unless additional assumptions are introduced (as typically done in the form of contractivity or monotonicity), solving MFG can in general be as hard as finding N -player general sum Nash.

We will prove that the problems are PPAD-complete, where PPAD is a class of computational problems studied in the seminal work by Papadimitriou (1994), containing the complete problem of finding N -player Nash equilibrium in general sum normal form games and finding the fixed point of continuous maps (Daskalakis et al., 2009; Chen et al., 2009). The class PPAD is conjectured to contain difficult problems with no polynomial time algorithms (Beame et al., 1995; Goldberg, 2011), hence our results can be seen as a proof of difficulty. Our results are significant since they imply that the MFG problems studied in literature are in

the same complexity class as general-sum N -player normal form games or N -player Markov games (Daskalakis et al., 2023). Once again, several computation-intensive aspects of our proofs will be postponed to the supplementary material.

Due to a technicality, we will prove the complexity results for a subset of possible reward and transition probability functions. We formalize this subset of possible rewards and dynamics as “simple” rewards/dynamics and also linear rewards, defined below.

Definition 4.1 (Simple/Linear Dynamics and Rewards). $R \in \mathcal{R}$ and $P \in \mathcal{P}$ are said to be *simple* if for any $s, s' \in \mathcal{S}$, $a \in \mathcal{A}$, $P(s'|s, a, \mu)$ and $R(s, a, \mu)$ are functions of μ that are expressible as finite combinations of arithmetic operations $+$, $-$, \times , \div and functions $\max\{\cdot, \cdot\}$, $\min\{\cdot, \cdot\}$ of coordinates of μ . They are called *linear* if $P(s'|s, a, \mu)$ and $R(s, a, \mu)$ are linear functions of μ for all s, a, s' . The set of simple rewards and dynamics are denoted by \mathcal{R}^{Sim} and \mathcal{P}^{Sim} respectively, and the set of linear rewards and transitions are denoted \mathcal{R}^{Lin} , \mathcal{P}^{Lin} respectively.

A note on simple functions. We define simple functions as above as in general there is no known efficient encoding of a Lipschitz continuous function as a sequence of bits. This is significant since a Turing machine accepts a finite sequence of bits as input. To solve this issue, we prove a slightly stronger hardness result that even games where $P(s'|s, a, \mu)$, $R(s, a, \mu)$ are Lipschitz functions with strong structure are PPAD-complete. Other larger classes of P, R including \mathcal{P}^{Sim} , \mathcal{R}^{Sim} will have similar intractability. See also arithmetic circuits with \max, \min gates (Daskalakis & Papadimitriou, 2011) for a similar idea.

4.1. The Complexity Class PPAD

The PPAD class is defined by the complete problem END-OF-THE-LINE (Daskalakis et al., 2009), whose formal definition we defer to the appendix as it is not used in proofs.

Definition 4.2 (PPAD, PPAD-hard, PPAD-complete). The class PPAD is defined as all search problems that can be reduced to END-OF-THE-LINE in polynomial time. If END-OF-THE-LINE can be reduced to a search problem \mathcal{S} in polynomial time, then \mathcal{S} is called PPAD-hard. A search problem \mathcal{S} is called PPAD-complete if it is both a member of PPAD and it is PPAD-hard.

While END-OF-THE-LINE defines the problem class PPAD, it is hard to construct direct reductions to it. We will instead use two problems that are known to be PPAD-complete (and hence can be equivalently used to define PPAD): solving generalized circuits and finding a NE for an N -player general sum game.

Definition 4.3 (Generalized Circuits (Rubinstein, 2015)). A generalized circuit $\mathcal{C} = (\mathcal{V}, \mathcal{G})$ is a finite set of nodes \mathcal{V} and gates \mathcal{G} . Each gate $G \in \mathcal{G}$ is characterized by the tuple

$G(\theta|v_1, v_2|v)$ where $G \in \{G_{\leftarrow}, G_{\times,+}, G_{<}\}$, $\theta \in \mathbb{R}^*$ is a parameter (possibly of length 0), $v_1, v_2 \in V \cup \{\perp\}$ are the input nodes (with \perp indicating an empty input) and $v \in V$ is the output node of the gate. The collection \mathcal{G} satisfies the property that if $G_1(\theta|v_1, v_2|v)$, $G_2(\theta'|v'_1, v'_2|v')$ $\in \mathcal{G}$ are distinct, then $v \neq v'$.

Such circuits define a set of constraints on values assigned to each gate, and finding such an assignment will be the associated computational problem for such a circuit description. We formally define the ε -GCIRCUIT problem to this end. ε -GCIRCUIT is a standard complete problem for the class PPAD, and we will work with it for our reductions. We will use the shorthand notation $x = y \pm \varepsilon$ to indicate that $x \in [y - \varepsilon, y + \varepsilon]$ for $x, y \in \mathbb{R}$.

Definition 4.4 (ε -GCIRCUIT (Rubinstein, 2015)). Given a generalized circuit $\mathcal{C} = (\mathcal{V}, \mathcal{G})$, a function $p : V \rightarrow [0, 1]$ is called an ε -satisfying assignment if:

- For every gate $G \in \mathcal{G}$ of the form $G_{\leftarrow}(\zeta|v)$ for $\zeta \in [0, 1]$, it holds that $p(v) = \zeta \pm \varepsilon$,
- For every gate $G \in \mathcal{G}$ of the form $G_{\times,+}(\alpha, \beta|v_1, v_2|v)$ for $\alpha, \beta \in [-1, 1]$, it holds that

$$p(v) \in [\max\{\min\{0, \alpha p(v_1) + \beta p(v_2)\}\}] \pm \varepsilon,$$

- For every gate $G \in \mathcal{G}$ of the form $G_{<}(|v_1, v_1|v)$ it holds that

$$p(v) = \begin{cases} 1 \pm \varepsilon, & p(v_1) \leq p(v_2) - \varepsilon, \\ 0 \pm \varepsilon, & p(v_1) \geq p(v_2) + \varepsilon. \end{cases}$$

The ε -GCIRCUIT problem is defined as follows:

*Given generalized circuit \mathcal{C} ,
find an ε -satisfying assignment of \mathcal{C} .*

ε -GCIRCUIT is one of the prototypical hard instances of PPAD problems as the result below suggests.

Theorem 4.5. (Rubinstein, 2015) *There exists $\varepsilon > 0$ such that ε -GCIRCUIT is PPAD-complete.*

In other words, ε -GCIRCUIT is representative of the most difficult problem in PPAD which suggests intractability. The ε -GCIRCUIT computational problem will be used in our proofs by reducing an arbitrary generalized circuit into solving a particular MFG.

We also use the general sum 2-player Nash computation problem, which is the standard problem of finding an approximate Nash equilibrium of a general sum bimatrix game.

Definition 4.6 (2-NASH). Given $\varepsilon > 0$, $K_1, K_2 \in \mathbb{N}_{>0}$, payoff matrices $A, B \in [0, 1]^{K_1, K_2}$, find an approximate Nash equilibrium $(\sigma_1, \sigma_2) \in \Delta_{K_1} \times \Delta_{K_2}$ such that

$$\max_{\sigma \in \Delta_{K_1}} U_A(\sigma, \sigma_2) - U_A(\sigma_1, \sigma_2) \leq \varepsilon,$$

$$\max_{\sigma \in \Delta_{K_2}} U_B(\sigma_1, \sigma) - U_B(\sigma_1, \sigma_2) \leq \varepsilon,$$

where $U_M(\sigma_1, \sigma_1) := \sum_{i \in [K_1]} \sum_{j \in [K_2]} M_{i,j} \sigma_1(i) \sigma_2(j)$ for any matrix $M \in [0, 1]^{K_1, K_2}$.

The following is the well-known result that even the 2-Nash general sum problem is PPAD-complete. In fact, any N -player general sum normal form game is PPAD-complete.

Theorem 4.7. (Chen et al., 2009) 2-NASH is PPAD-complete.

4.2. Complexity of Stat-MFG

Next, we provide our hardness results for the Stat-MFG problem. Notably, for Stat-MFG, the stability subproblem of finding a stable distribution for a fixed policy π itself is PPAD-hard. Even without considering the optimality conditions, finding a stable distribution in general for a fixed policy is intractable, without additional assumptions (e.g. Γ_P is contractive or non-expansive). We define the computational problem below and state the results.

Definition 4.8 (ε -STATDIST). Given finite state-action sets \mathcal{S}, \mathcal{A} , simple dynamics $P \in \mathcal{P}^{\text{Sim}}$ and policy π , find $\mu^* \in \Delta_{\mathcal{S}}$ such that $\|\Gamma_P(\mu^*, \pi) - \mu^*\|_{\infty} \leq \frac{\varepsilon}{|\mathcal{S}|}$.

The computational problem as described above is to find an approximate fixed point of $\Gamma_P(\cdot, \pi)$ which corresponds to an approximate stable distribution of policy π . We show that ε -STATDIST is PPAD-complete for some fixed constant ε .

Theorem 4.9 (ε -STATDIST is PPAD-complete). For some $\varepsilon > 0$, the problem ε -STATDIST is PPAD-complete.

Consequently, there is no polynomial time algorithm for ε -STATDIST unless PPAD=P, which is conjectured to be not the case.

Corollary 4.10. There exists a $\varepsilon > 0$ such that there exists no polynomial time algorithm for ε -STATDIST, unless P = PPAD.

Most notably, these results show that the stable distribution oracle of (Cui & Koepl, 2021) might be intractable to compute in general, and the shared assumption that $\Gamma_P(\cdot, \pi)$ is contractive in some norm found in many works (Xie et al., 2021; Anahtarci et al., 2022; Yardim et al., 2023a) might not be trivial to remove without sacrificing tractability.

4.3. Complexity of FH-MFG

We will show that finding an ε solution to the finite horizon problem is also PPAD-complete, in particular even if we

restrict our attention to the case when $H = 2$ and the transition probabilities P do not depend on μ . We formalize the structured computational FH-MFG problem.

Definition 4.11 ((ε, H) -FH-NASH). Given simple reward function $R \in \mathcal{R}^{\text{Sim}}$, transition matrix $P(s'|s, a)$, and initial distribution $\mu_0 \in \Delta_{\mathcal{S}}$, find a time dependent policy $\{\pi_h\}_{h=0}^{H-1}$ such that $\mathcal{E}_{P,R}^H(\{\pi_h\}_{h=0}^{H-1}) \leq \varepsilon/|\mathcal{S}|$.

Our main result for FH-MFG is that even in the case of $H = 2$, the problem is PPAD-complete.

Theorem 4.12 ($(\varepsilon, 2)$ -FH-NASH is PPAD-complete). There exists an $\varepsilon > 0$ such that the problem $(\varepsilon, 2)$ -FH-NASH is PPAD-complete.

Corollary 4.13. There exists a $\varepsilon > 0$ such that there exists no polynomial time algorithm for $(\varepsilon, 2)$ -FH-NASH, unless P = PPAD.

These results for the FH-MFG show that the (weak) monotonicity assumption present in works such as (Perrin et al., 2020; Pérolat et al., 2022) might also be necessary, as in the absence of any structural assumptions the problems are provably hard.

Finally, we also show that even if $R(s, a, \mu)$ is a linear function of μ for all s, a (that is, $R \in \mathcal{R}^{\text{Lin}}$), the intractability holds, although not for fixed ε . This follows from a reduction to 2-NASH. We define the linear computational problem below.

Definition 4.14 (H -FH-LINEAR). Given $\varepsilon > 0$, linear reward function $R \in \mathcal{R}^{\text{Lin}}$, transition matrix $P(s'|s, a)$, find a time dependent policy $\{\pi_h\}_{h=0}^{H-1}$ such that $\mathcal{E}_{P,R}^H(\{\pi_h\}_{h=0}^{H-1}) \leq \varepsilon$.

Theorem 4.15 (2-FH-LINEAR is PPAD-complete). The problem 2-FH-LINEAR is PPAD-complete.

We emphasize that for 2-FH-LINEAR the accuracy ε is also an input of the problem: hence the existence of a pseudo-polynomial time algorithm is not ruled out.

5. Discussion and Conclusion

We provided novel results on when mean-field RL is relevant for real-world applications and when it is tractable from a computational perspective. Our results differ from existing work by provably characterizing cases where MFGs might have practical shortcomings. From the approximation perspective, we show clear conditions and lower bounds on when the MFGs efficiently approximate real-world games. Computationally, we show that even simple MFGs can be as hard as solving N -player general sum games.

We emphasize that our results do not discard MFGs, but rather identify potential bottlenecks (and conditions to overcome these) when using mean-field RL to compute a good approximate NE.

References

- 440
441
442 Anahtarci, B., Kariksiz, C. D., and Saldi, N. Q-learning
443 in regularized mean-field games. *Dynamic Games and*
444 *Applications*, pp. 1–29, 2022.
- 445 Beame, P., Cook, S., Edmonds, J., Impagliazzo, R., and
446 Pitassi, T. The relative complexity of np search problems.
447 In *Proceedings of the twenty-seventh annual ACM sym-*
448 *posium on Theory of computing*, pp. 303–314, Las Vegas,
449 Nevada, USA, 1995.
- 450
451 Carmona, R. and Delarue, F. Probabilistic analysis of mean-
452 field games. *SIAM Journal on Control and Optimization*,
453 51(4):2705–2734, 2013.
- 454 Carmona, R., Delarue, F., et al. *Probabilistic theory of mean*
455 *field games with applications I-II*. Springer, 2018.
- 456
457 Chen, X., Deng, X., and Teng, S.-H. Settling the complexity
458 of computing two-player nash equilibria. *Journal of the*
459 *ACM (JACM)*, 56(3):1–57, 2009.
- 460
461 Cui, K. and Koepl, H. Approximately solving mean field
462 games via entropy-regularized deep reinforcement learn-
463 ing. In *International Conference on Artificial Intelligence*
464 *and Statistics*, pp. 1909–1917. PMLR, 2021.
- 465
466 Daskalakis, C. and Papadimitriou, C. Continuous local
467 search. In *Proceedings of the twenty-second annual ACM-*
468 *SIAM symposium on Discrete Algorithms*, pp. 790–804.
469 SIAM, 2011.
- 470
471 Daskalakis, C., Goldberg, P. W., and Papadimitriou, C. H.
472 The complexity of computing a nash equilibrium. *Com-*
473 *munications of the ACM*, 52(2):89–97, 2009.
- 474
475 Daskalakis, C., Golowich, N., and Zhang, K. The complex-
476 ity of markov equilibrium in stochastic games. In *The*
477 *Thirty Sixth Annual Conference on Learning Theory*, pp.
478 4180–4234. PMLR, 2023.
- 479 Geist, M., Pérolat, J., Laurière, M., Elie, R., Perrin, S.,
480 Bachem, O., Munos, R., and Pietquin, O. Concave utility
481 reinforcement learning: The mean-field game viewpoint.
482 In *Proceedings of the 21st International Conference on*
483 *Autonomous Agents and Multiagent Systems, AAMAS*
484 *'22*, pp. 489–497, Richland, SC, 2022. International Founda-
485 tion for Autonomous Agents and Multiagent Systems. ISBN
486 9781450392136.
- 487
488 Goldberg, P. W. A survey of ppad-completeness for com-
489 puting nash equilibria. *arXiv preprint arXiv:1103.2709*,
490 2011.
- 491
492 Guo, X., Hu, A., Xu, R., and Zhang, J. Learning mean-
493 field games. *Advances in Neural Information Processing*
494 *Systems*, 32, 2019.
- Guo, X., Hu, A., Xu, R., and Zhang, J. A general frame-
work for learning mean-field games. *Mathematics of*
Operations Research, 2022a.
- Guo, X., Hu, A., and Zhang, J. Mf-omo: An optimiza-
tion formulation of mean-field games. *arXiv preprint*
arXiv:2206.09608, 2022b.
- Huang, J., Yardim, B., and He, N. On the statisti-
cal efficiency of mean field reinforcement learning
with general function approximation. *arXiv preprint*
arXiv:2305.11283, 2023.
- Huang, M., Malhamé, R. P., and Caines, P. E. Large pop-
ulation stochastic dynamic games: closed-loop mckean-
vaslov systems and the nash certainty equivalence prin-
ciple. *Communications in Information & Systems*, 6(3):
221–252, 2006.
- Iyer, K., Johari, R., and Sundararajan, M. Mean field equi-
libria of dynamic auctions with learning. *Management*
Science, 60(12):2949–2970, 2014.
- Kaas, R. and Buhman, J. M. Mean, median and mode
in binomial distributions. *Statistica Neerlandica*, 34(1):
13–18, 1980.
- Lasry, J.-M. and Lions, P.-L. Mean field games. *Japanese*
journal of mathematics, 2(1):229–260, 2007.
- Laurière, M., Perrin, S., Girgin, S., Muller, P., Jain, A., Ca-
bannes, T., Piliouras, G., P’erolat, J., Elie, R., Pietquin,
O., and Geist, M. Scalable deep reinforcement learn-
ing algorithms for mean field games. In *International*
Conference on Machine Learning, 2022.
- Mao, W., Qiu, H., Wang, C., Franke, H., Kalbarczyk, Z.,
Iyer, R., and Basar, T. A mean-field game approach to
cloud resource management with function approximation.
In *Advances in Neural Information Processing Systems*,
2022.
- Matignon, L., Laurent, G. J., and Le Fort-Piat, N. Hys-
teretic q-learning: an algorithm for decentralized rein-
forcement learning in cooperative multi-agent teams. In
2007 IEEE/RSJ International Conference on Intelligent
Robots and Systems, pp. 64–69. IEEE, 2007.
- McDiarmid, C. et al. On the method of bounded differences.
Surveys in combinatorics, 141(1):148–188, 1989.
- Papadimitriou, C. H. On the complexity of the parity argu-
ment and other inefficient proofs of existence. *Journal of*
Computer and system Sciences, 48(3):498–532, 1994.
- Perolat, J., Scherrer, B., Piot, B., and Pietquin, O. Approx-
imate dynamic programming for two-player zero-sum
markov games. In *International Conference on Machine*
Learning, pp. 1321–1329. PMLR, 2015.

- 495 Pérolat, J., Perrin, S., Elie, R., Laurière, M., Piliouras, G.,
496 Geist, M., Tuyls, K., and Pietquin, O. Scaling mean field
497 games by online mirror descent. In *Proceedings of the*
498 *21st International Conference on Autonomous Agents and*
499 *Multiagent Systems*, pp. 1028–1037, 2022.
- 500
501 Perrin, S., Pérolat, J., Laurière, M., Geist, M., Elie, R.,
502 and Pietquin, O. Fictitious play for mean field games:
503 Continuous time analysis and applications. *Advances*
504 *in Neural Information Processing Systems*, 33:13199–
505 13213, 2020.
- 506
507 Rashedi, N., Tajeddini, M. A., and Kebriaei, H. Markov
508 game approach for multi-agent competitive bidding strate-
509 gies in electricity market. *IET Generation, Transmission*
510 *& Distribution*, 10(15):3756–3763, 2016.
- 511
512 Rubinstein, A. Inapproximability of nash equilibrium. In
513 *Proceedings of the forty-seventh annual ACM symposium*
514 *on Theory of computing*, pp. 409–418, 2015.
- 515
516 Saldi, N., Basar, T., and Raginsky, M. Markov–nash equi-
517 libria in mean-field games with discounted cost. *SIAM*
518 *Journal on Control and Optimization*, 56(6):4256–4287,
519 2018.
- 520
521 Samvelyan, M., Rashid, T., Schroeder de Witt, C., Farquhar,
522 G., Nardelli, N., Rudner, T. G. J., Hung, C.-M., Torr,
523 P. H. S., Foerster, J., and Whiteson, S. The starcraft
524 multi-agent challenge. In *Proc. of the 18th International*
525 *Conference on Autonomous Agents and Multiagent Sys-*
526 *tems (AAMAS 2019)*, AAMAS ’19, pp. 2186–2188, Rich-
527 land, SC, 2019. International Foundation for Autonomous
528 Agents and Multiagent Systems.
- 529
530 Shavandi, A. and Khedmati, M. A multi-agent deep rein-
531 forcement learning framework for algorithmic trading in
532 financial markets. *Expert Systems with Applications*, 208:
533 118124, 2022.
- 534
535 Wang, L., Yang, Z., and Wang, Z. Breaking the curse of
536 many agents: Provable mean embedding q-iteration for
537 mean-field reinforcement learning. In *International con-*
538 *ference on machine learning*, pp. 10092–10103. PMLR,
539 2020.
- 540
541 Wiering, M. A. Multi-agent reinforcement learning for
542 traffic light control. In *Machine Learning: Proceedings of*
543 *the Seventeenth International Conference (ICML’2000)*,
544 pp. 1151–1158, 2000.
- 545
546 Xie, Q., Yang, Z., Wang, Z., and Minca, A. Learning while
547 playing in mean-field games: Convergence and optimality.
548 In *International Conference on Machine Learning*, pp.
549 11436–11447. PMLR, 2021.
- Yardim, B., Cayci, S., Geist, M., and He, N. Policy mirror
ascent for efficient and independent learning in mean
field games. In *International Conference on Machine*
Learning, pp. 39722–39754. PMLR, 2023a.
- Yardim, B., Cayci, S., and He, N. Stateless mean-field
games: A framework for independent learning with large
populations. In *Sixteenth European Workshop on Rein-*
forcement Learning, 2023b.
- Zaman, M. A. U., Koppel, A., Bhatt, S., and Basar, T.
Oracle-free reinforcement learning in mean-field games
along a single sample path. In *International Conference*
on Artificial Intelligence and Statistics, pp. 10178–10206.
PMLR, 2023.

A. MFG Approximation Results

A.1. Preliminaries

To establish explicit upper bounds on the approximation rate, we will use standard concentration tools.

Definition A.1 (Sub-Gaussian). Random variable ξ is called sub-Gaussian with variance proxy σ^2 if $\forall \lambda \in \mathbb{R} : \mathbb{E} [e^{\lambda(\xi - \mathbb{E}[\xi])}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}$. In this case, we write $\xi \in SG(\sigma^2)$.

It is easy to show that if $\xi \in SG(\sigma^2)$, then $\alpha\xi \in SG(\alpha^2\sigma^2)$ for any constant $\alpha \in \mathbb{R}$. Furthermore, if ξ_1, \dots, ξ_n are independent random variables with $\xi_i \in SG(\sigma_i^2)$, then $\sum_i \xi_i \in SG(\sum_i \sigma_i^2)$. Finally, if ξ is almost surely bounded in $[a, b]$, then $\xi_i \in SG((b-a)^2/4)$. We also state the well-known Hoeffding concentration bound and a corollary, Lemma A.3.

Lemma A.2 (Hoeffding inequality (McDiarmid et al., 1989)). *Let $\xi \in SG(\sigma^2)$. Then for any $t > 0$ it holds that $\mathbb{P}(|\xi - \mathbb{E}[\xi]| \geq t) \leq 2e^{-\frac{t^2}{2\sigma^2}}$.*

Lemma A.3. *Let $\xi \in SG(\sigma^2)$. Then*

$$\mathbb{E}[|\xi - \mathbb{E}[\xi]|] \leq \sqrt{2\pi\sigma^2}, \quad \mathbb{E}[(\xi - \mathbb{E}[\xi])^2] \leq 4\sigma^2$$

Proof.

$$\begin{aligned} \mathbb{E}[|\xi - \mathbb{E}[\xi]|] &= \int_0^\infty \mathbb{P}(|\xi - \mathbb{E}[\xi]| \geq t) dt \\ &\stackrel{(I)}{\leq} 2 \int_0^\infty e^{-\frac{t^2}{2\sigma^2}} dt = \sqrt{2\pi\sigma^2} \end{aligned}$$

Inequality (I) is true due to Lemma A.2. Likewise,

$$\begin{aligned} \mathbb{E}[(\xi - \mathbb{E}[\xi])^2] &= \int_0^\infty \mathbb{P}((\xi - \mathbb{E}[\xi])^2 \geq t) dt \\ &= \int_0^\infty \mathbb{P}(|\xi - \mathbb{E}[\xi]| \geq \sqrt{t}) dt \\ &\stackrel{(II)}{\leq} 2 \int_0^\infty e^{-\frac{t}{2\sigma^2}} dt = 4\sigma^2 \end{aligned}$$

□

Establishing lower bounds for the mean-field approximation of the N -player game will be more challenging as it will require different tools. To establish lower bounds, we will need to use the following anti-concentration result for the binomial distribution.

Lemma A.4 (Anti-concentration for binomial). *Let $N \in \mathbb{N}_{>0}$ and $X \sim \text{Binom}(N, p)$ be drawn from a binomial distribution for some $p \in [1/2, 1]$. Then, $\mathbb{P}\left[X \geq \frac{N}{2} + \frac{\sqrt{N}}{2}\right] \geq \frac{1}{20}$.*

Proof. For $k_0 := \left\lceil \frac{N}{2} + \frac{\sqrt{N}}{2} \right\rceil$, we will lower bound $\sum_{k=k_0}^N \binom{N}{k} p^k (1-p)^{N-k}$ when N is large enough. If $k_0 < \lceil Np \rceil$, then the probability in the statement above is bounded below trivially by $1/2$ since $\lceil Np \rceil$ lower bounds the median of the binomial (Kaas & Buhrman, 1980). Otherwise, if $k_0 \geq \lceil Np \rceil$, then the function $\bar{p} \rightarrow \bar{p}^k (1-\bar{p})^{N-k}$ is increasing in \bar{p} in the interval $[0, p]$. As $1/2 \in [0, p]$, it is then sufficient to assume $p = 1/2$, and to upper bound $\mathbb{P}\left[\frac{N}{2} - \frac{\sqrt{N}}{2} < X < \frac{N}{2} + \frac{\sqrt{N}}{2}\right]$ by $9/10$ as the binomial probability mass is symmetric around $\frac{N}{2}$ when $p = 1/2$.

First assuming N is even, we obtain by monotonicity $\binom{N}{k} \leq \binom{N}{N/2}$. Using the Stirling bound $\sqrt{2\pi} k^{k+\frac{1}{2}} e^{-k} \leq k! \leq e k^{k+\frac{1}{2}} e^{-k}$, we further upper bound $\binom{N}{N/2} \leq \frac{e}{\pi} \frac{2^N}{\sqrt{N}}$, resulting in the bound $\mathbb{P}\left[\frac{N}{2} - \frac{\sqrt{N}}{2} < X < \frac{N}{2} + \frac{\sqrt{N}}{2}\right] \leq 2^{-N} \sqrt{N} \binom{N}{N/2} \leq \frac{e}{\pi} \leq 9/10$, since there are at most \sqrt{N} binomial coefficients being summed. Finally, assume $N = 2m + 1$ is odd, then by the binomial formula $\binom{2m+1}{m+1} = \binom{2m}{m+1} + \binom{2m}{m} \leq 2 \binom{2m}{m} \leq \frac{2e}{\pi} \frac{2^{2m}}{\sqrt{2m}}$. Hence we have the bound on the sum

$\mathbb{P}\left[\frac{N}{2} - \frac{\sqrt{N}}{2} < X < \frac{N}{2} + \frac{\sqrt{N}}{2}\right] \leq \frac{e\sqrt{N}}{\pi} \frac{1}{\sqrt{N-1}}$. It is easy to verify that for $N \geq 16$, $\frac{e\sqrt{N}}{\pi\sqrt{N-1}} \leq 9/10$, and the case when $N < 16$ and N is odd follows by manual computation. \square

Finally, we prove slightly more general upper bounds than presented in the main text that approximates the exploitability of an *approximate* MFG-NE in a finite population setting. Hence we define the following notions approximate FH-MFG and Stat-MFG.

Definition A.5 (δ -FH-MFG-NE). Let $(\mathcal{S}, \mathcal{A}, H, P, R, \mu_0)$ be a FH-MFG. Then, a δ -FH-MFG Nash equilibrium is defined as:

$$\begin{aligned}
 \text{Policy } \pi_\delta^* &= \{\pi_{\delta,h}^*\}_{h=0}^{H-1} \in \Pi_H \text{ such that} \\
 \mathcal{E}_{P,R}^H(\{\pi_{\delta,h}^*\}_{h=0}^{H-1}) &\leq \delta. \tag{\delta\text{-FH-MFG-NE}}
 \end{aligned}$$

Definition A.6 (δ -Stat-MFG-NE). Let $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ be a Stat-MFG. A policy-population pair $(\mu_\delta^*, \pi_\delta^*) \in \Delta_{\mathcal{S}} \times \Pi$ is called a δ -Stat-MFG Nash equilibrium if the two conditions hold:

$$\begin{aligned}
 \text{Stability: } \mu_\delta^* &= \Gamma_P(\mu_\delta^*, \pi_\delta^*), \\
 \text{Optimality: } V_{P,R}^\gamma(\mu_\delta^*, \pi_\delta^*) &\geq \max_{\pi \in \Pi} V_{P,R}^\gamma(\mu_\delta^*, \pi) - \delta. \tag{\delta\text{-Stat-MFG-NE}}
 \end{aligned}$$

A.2. Upper Bound for FH-MFG: Extended Proof of Theorem 3.2

Throughout this section we work with fixed $P \in \mathcal{P}_{K_\mu}$ and $R \in \mathcal{R}_{L_\mu}$. For any \mathcal{X} valued random variable x denote $\mathcal{L}(x)(\cdot) \in \Delta_{\mathcal{X}}$ as the distribution of x . We start by introducing some notation.

For given R and P define the following constants:

$$\begin{aligned}
 L_s &:= \sup_{s,s',a,\mu} |R(s,a,\mu) - R(s',a,\mu)|, \\
 L_a &:= \sup_{s,a,a',\mu} |R(s,a,\mu) - R(s,a',\mu)|, \\
 K_s &:= \sup_{s,s',a,\mu} \|P(\cdot|s,a,\mu) - P(\cdot|s',a,\mu)\|, \\
 K_a &:= \sup_{s,a,a',\mu} \|P(\cdot|s,a,\mu) - P(\cdot|s,a',\mu)\|.
 \end{aligned}$$

R and P are bounded due to Definition 2.1, thus all constants K_a, K_s, L_a, L_s are finite and well-defined, and it always holds that $K_s, K_a \leq 2$ and $L_s, L_a \leq 1$. With the above definition of constants, the more general Lipschitz condition holds: $\forall s, s' \in \mathcal{S}, a, a' \in \mathcal{A}, \mu, \mu' \in \Delta_{\mathcal{S}}$

$$\begin{aligned}
 \|P(\cdot|s,a,\mu) - P(\cdot|s',a',\mu')\|_1 &\leq K_\mu \|\mu - \mu'\|_1 + K_s d(s, s') \\
 &\quad + K_a d(a, a'), \\
 |R(s,a,\mu) - R(s',a',\mu')| &\leq L_\mu \|\mu - \mu'\|_1 + L_s d(s, s') \\
 &\quad + L_a d(a, a').
 \end{aligned}$$

We also introduce the shorthand notation for any $s \in \mathcal{S}, u \in \Delta_{\mathcal{A}}, \mu \in \Delta_{\mathcal{S}}$:

$$\begin{aligned}
 \bar{P}(\cdot|s, u, \mu) &:= \sum_{a \in \mathcal{A}} u(a) P(\cdot|s, a, \mu), \\
 \bar{R}(s, u, \mu) &:= \sum_{a \in \mathcal{A}} u(a) R(s, a, \mu).
 \end{aligned}$$

By (?)Lemma C.1]yardim2023policy, it holds that

$$\begin{aligned}
 \|\bar{P}(\cdot|s, u, \mu) - \bar{P}(\cdot|s', u', \mu')\|_1 &\leq K_\mu \|\mu - \mu'\|_1 + K_s d(s, s') \\
 &\quad + \frac{K_a}{2} \|u - u'\|_1, \\
 |\bar{R}(s, u, \mu) - \bar{R}(s', u', \mu')| &\leq L_\mu \|\mu - \mu'\|_1 + L_s d(s, s') \\
 &\quad + \frac{L_a}{2} \|u - u'\|_1. \tag{1}
 \end{aligned}$$

We will define a new operator for tracking the evolution of the population distribution over finite time horizons for a time-varying policy $\forall \boldsymbol{\pi} = \{\pi_h\}_{h=0}^{H-1} \in \Pi_H$:

$$\begin{aligned}\Gamma_P^h(\mu, \boldsymbol{\pi}) &:= \underbrace{\Gamma_P(\dots \Gamma_P(\Gamma_P(\mu, \pi_0), \pi_1) \dots, \pi_{h-1})}_{h \text{ times}} \\ &= \mu_h^\pi = \Lambda_P^H(\mu_0, \boldsymbol{\pi})_h,\end{aligned}$$

so $\Gamma_P^0(\mu, \boldsymbol{\pi}) = \mu_0$. By repeated applications of Lemma 2.2, we obtain the Lipschitz condition:

$$\begin{aligned}\|\Gamma_P^n(\mu, \{\pi_i\}_{i=0}^{n-1}) - \Gamma_P^n(\mu', \{\pi'_i\}_{i=0}^{n-1})\|_1 &\leq L_{pop, \mu} \|\Gamma_P^{n-1}(\mu, \{\pi_i\}_{i=0}^{n-2}) - \Gamma_P^{n-1}(\mu', \{\pi'_i\}_{i=0}^{n-2})\|_1 \\ &\quad + \frac{K_a}{2} \|\pi_{n-1} - \pi'_{n-1}\|_1 \\ &\leq L_{pop, \mu}^n \|\mu - \mu'\|_1 + \frac{K_a}{2} \sum_{i=0}^{n-1} L_{pop, \mu}^{n-1-i} \|\pi_i - \pi'_i\|_1,\end{aligned}\tag{2}$$

where $L_{pop, \mu} = (K_\mu + \frac{K_s}{2} + \frac{K_a}{2})$.

The proof will proceed in three steps:

- **Step 1.** Bounding the expected deviation of the empirical population distribution from the mean-field distribution $\mathbb{E}[\|\hat{\mu}_h - \mu_h^\pi\|_1]$ for any given policy $\boldsymbol{\pi}$.
- **Step 2.** Bounding difference of N agent value function $J_{P,R}^{H,N,(i)}$ and the infinite player value function $V_{P,R}^H$.
- **Step 3.** Bounding the exploitability of an agent when each of N agents are playing the FH-MFG-NE policy.

Step 1: Empirical distribution bound. Due to its relevance for a general connection between the FH-MFG and the N -player game, we state this result in the form of an explicit bound.

Lemma A.7. *Suppose for the N -FH-MFG $(N, \mathcal{S}, \mathcal{A}, N, P, R, \gamma)$, agents $i = 1, \dots, N$ follow policies $\boldsymbol{\pi}^i = \{\pi_h^i\}_h$. Let $\bar{\boldsymbol{\pi}} = \{\bar{\pi}_h\}_h \in \Pi^H$ be arbitrary and $\mu^{\bar{\boldsymbol{\pi}}} := \{\mu_h^{\bar{\boldsymbol{\pi}}}\}_{h=0}^{H-1} = \Lambda_P^H(\mu_0, \bar{\boldsymbol{\pi}})$. Then for all $h \in \{0, \dots, H-1\}$, it holds that:*

$$\mathbb{E}[\|\hat{\mu}_h - \mu_h^{\bar{\boldsymbol{\pi}}}\|_1] \leq \frac{1 - L_{pop, \mu}^{h+1}}{1 - L_{pop, \mu}} |\mathcal{S}| \sqrt{\frac{\pi}{2N}} + \frac{K_a}{2N} \sum_{i=0}^{h-1} L_{pop, \mu}^{h-i-1} \Delta_{\pi_i},$$

where $\Delta_h := \frac{1}{N} \sum_i \|\bar{\pi}_h - \pi_h^i\|_1$

Proof. The proof will proceed inductively over h . First, for time $h = 0$, we have

$$\mathbb{E}[\|\hat{\mu}_0 - \mu_0\|_1] = \sum_{s \in \mathcal{S}} \mathbb{E} \left[\left| \frac{1}{N} \sum_{i=1}^N (\mathbb{1}_{\{s_0^i=s\}} - \mu_0(s)) \right| \right] \leq |\mathcal{S}| \sqrt{\frac{\pi}{2N}},$$

where the last line is due to Lemma A.3 and the fact that $\mathbb{1}_{\{s_0^i=s\}}$ are bounded (hence subgaussian) random variables, and that in the finite state space we have $\mathbb{E}[\mathbb{1}_{\{s_0^i=s\}}] = \mu_0(s)$.

Next, denoting the σ -algebra induced by the random variables $(\{s_h^i\})_{i, h' \leq h}$ as \mathcal{F}_h , we have that:

$$\begin{aligned}&\mathbb{E}[\|\hat{\mu}_{h+1} - \mu_{h+1}^{\bar{\boldsymbol{\pi}}}\|_1 | \mathcal{F}_h] \\ &\leq \underbrace{\mathbb{E}[\|\mathbb{E}[\hat{\mu}_{h+1} | \mathcal{F}_h] - \Gamma_P(\hat{\mu}_h, \bar{\boldsymbol{\pi}}_h)\|_1 | \mathcal{F}_h]}_{(\square)} \\ &\quad + \underbrace{\mathbb{E}[\|\hat{\mu}_{h+1} - \mathbb{E}[\hat{\mu}_{h+1} | \mathcal{F}_h]\|_1 | \mathcal{F}_h]}_{(\triangle)} + \underbrace{\mathbb{E}[\|\Gamma_P(\hat{\mu}_h, \bar{\boldsymbol{\pi}}_h) - \mu_{h+1}^{\bar{\boldsymbol{\pi}}}\|_1 | \mathcal{F}_h]}_{(\heartsuit)}\end{aligned}\tag{3}$$

We upper bound the three terms separately. For (Δ) , it holds that

$$\begin{aligned} (\Delta) &= \mathbb{E} [\|\widehat{\mu}_{h+1} - \mathbb{E}[\widehat{\mu}_{h+1} | \mathcal{F}_h]\|_1 | \mathcal{F}_h] \\ &= \sum_{s \in \mathcal{S}} \mathbb{E} [|\widehat{\mu}_{h+1}(s) - \mathbb{E}[\widehat{\mu}_{h+1}(s) | \mathcal{F}_h]| | \mathcal{F}_h] \leq |\mathcal{S}| \sqrt{\frac{\pi}{2N}}, \end{aligned}$$

since each $\widehat{\mu}_{h+1}(s)$ is an average of independent subgaussian random variables given \mathcal{F}_h . Specifically, each indicator is bounded $\mathbb{1}_{\{s_{h+1}^i=s\}} \in [0, 1]$ a.s. and therefore is sub-Gaussian with $\mathbb{1}_{\{s_{h+1}^i=s\}} \in SG(1/4)$. Thus we get $\widehat{\mu}_{h+1}(s) \in SG(1/(4N))$ and apply bound on expected value discussed in Appendix A.1.

Next, for $(\square) = \|\mathbb{E}[\widehat{\mu}_{h+1} | \mathcal{F}_h] - \Gamma_P(\widehat{\mu}_h, \bar{\pi}_h)\|_1$, we note that

$$\mathbb{E}[\widehat{\mu}_{h+1}(s) | \mathcal{F}_h] = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{s_{h+1}^i=s\}} | \mathcal{F}_h \right] = \frac{1}{N} \sum_{i=1}^N \bar{P}(s | s_h^i, \pi_h^i(s_h^i), \widehat{\mu}_h),$$

therefore

$$\begin{aligned} (\square) &= \left\| \frac{1}{N} \sum_{i=1}^N \bar{P}(\cdot | s_h^i, \pi_h^i(\cdot | s_h^i), \widehat{\mu}_h) - \sum_{s'} \widehat{\mu}_h(s') \bar{P}(\cdot | s', \pi_h(\cdot | s'), \widehat{\mu}_h) \right\|_1 \\ &= \left\| \frac{1}{N} \sum_{i=1}^N (\bar{P}(\cdot | s_h^i, \pi_h^i(\cdot | s_h^i), \widehat{\mu}_h) - \bar{P}(\cdot | s_h^i, \pi_h(\cdot | s_h^i), \widehat{\mu}_h)) \right\|_1 \\ &\leq \frac{1}{N} \sum_{i=1}^N \|\bar{P}(\cdot | s_h^i, \pi_h^i(\cdot | s_h^i), \widehat{\mu}_h) - \bar{P}(\cdot | s_h^i, \pi_h(\cdot | s_h^i), \widehat{\mu}_h)\|_1 \\ &\stackrel{(I)}{\leq} \frac{K_a}{2N} \sum_{i=1}^N \|\pi_h^i(\cdot | s_h^i) - \pi_h(\cdot | s_h^i)\|_1 \leq \frac{K_a}{2} \Delta_h, \end{aligned}$$

where (I) follows from the Lipschitz property (1). Finally, the last term (\heartsuit) can be bounded using:

$$(\heartsuit) = \mathbb{E} [\|\Gamma_P(\widehat{\mu}_h, \bar{\pi}_h) - \Gamma_P(\mu_h^{\bar{\pi}}, \bar{\pi}_h)\|_1 | \mathcal{F}_h] \leq L_{pop, \mu} \|\widehat{\mu}_h - \mu_h^{\bar{\pi}}\|_1.$$

To conclude, merging the bounds on the three terms in Inequality (3) and taking the expectations we obtain:

$$\mathbb{E} [\|\widehat{\mu}_{h+1} - \mu_{h+1}^{\bar{\pi}}\|_1] \leq L_{pop, \mu} \mathbb{E} [\|\widehat{\mu}_h - \mu_h^{\bar{\pi}}\|_1] + |\mathcal{S}| \sqrt{\frac{\pi}{2N}} + \frac{K_a \Delta_h}{2}.$$

Induction on h yields the statement of the lemma. \square

Step 2: Bounding difference of N agent value function. Next, we bound the difference between the N -player expected reward function $J_{P,R}^{H,N,(1)}$ and the infinite player expected reward function $V_{P,R}^H$. For ease of reading, expectations, probabilities, and laws of random variables will be denoted $\mathbb{E}_\infty, \mathbb{P}_\infty, \mathcal{L}_\infty$ respectively over the infinite player finite horizon game and $\mathbb{E}_N, \mathbb{P}_N, \mathcal{L}_N$ respectively over the N -player game. We use the regular notation $\mathbb{E}[\cdot], \mathbb{P}[\cdot], \mathcal{L}(\cdot)$ without subscripts if the underlying randomness is clearly defined. We state the main result of this step in the following lemma.

Lemma A.8. *Suppose N -FH-MFG agents follow the same sequence of policies $\boldsymbol{\pi} = \{\pi_h\}_{h=0}^{H-1}$. Then*

$$\begin{aligned} &\left| J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}, \dots, \boldsymbol{\pi}) - V_{P,R}^H(\Lambda_P^H(\mu_0, \boldsymbol{\pi}), \boldsymbol{\pi}) \right| \\ &\leq (L_\mu + \frac{L_s}{2}) |\mathcal{S}| \sqrt{\frac{\pi}{2N}} \sum_{h=0}^{H-1} \frac{1 - L_{pop, \mu}^{h+1}}{1 - L_{pop, \mu}}. \end{aligned}$$

Proof. Due to symmetry in the N agent game, any permutation $\sigma : [N] \rightarrow [N]$ of agents does not change their distribution, that is $\mathcal{L}_N(s_h^1, \dots, s_h^N) = \mathcal{L}_N(s_h^{\sigma(1)}, \dots, s_h^{\sigma(N)})$. We can then conclude that:

$$\begin{aligned} \mathbb{E}_N [R(s_h^1, a_h^1, \hat{\mu}_h)] &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}_N [R(s_h^i, a_h^i, \hat{\mu}_h)] \\ &= \mathbb{E}_N \left[\sum_{s \in \mathcal{S}} \hat{\mu}_h(s) \bar{R}(s, \pi_h(s), \hat{\mu}_h) \right] \end{aligned}$$

Therefore, we by definition:

$$J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}, \dots, \boldsymbol{\pi}) = \mathbb{E}_N \left[\sum_{h=0}^{H-1} \sum_{s \in \mathcal{S}} \hat{\mu}_h(s) \bar{R}(s, \pi_h(s), \hat{\mu}_h) \right].$$

Next, in the FH-MFG, under the population distribution $\{\mu_h\}_{h=0}^{H-1} = \Lambda_P^H(\mu_0, \boldsymbol{\pi})$ we have that for all $h \in 0, \dots, H-1$,

$$\begin{aligned} \mathbb{P}_\infty(s_0 = \cdot) &= \mu_0, \\ \mathbb{P}_\infty(s_{h+1} = \cdot) &= \sum_{s \in \mathcal{S}} \mathbb{P}_\infty(s_h = s) \mathbb{P}_\infty(s_h = \cdot | s_h = s) \\ &= \Gamma_P(\mathbb{P}_\infty(s_h = \cdot), \pi_h), \end{aligned}$$

so by induction $\mathbb{P}_\infty(s_h = \cdot) = \mu_h$. Then we can conclude that

$$\begin{aligned} V_{P,R}^H(\Lambda_P^H(\mu_0, \boldsymbol{\pi}), \boldsymbol{\pi}) &= \mathbb{E}_\infty \left[\sum_{h=0}^{H-1} R(s_h, \pi_h(s_h), \mu_h) \right] \\ &= \sum_{h=0}^{H-1} \sum_{s \in \mathcal{S}} \mu_h(s) R(s, \pi_h(s), \mu_h). \end{aligned}$$

Merging the two equalities for J, V , we have the bound:

$$\begin{aligned} &|J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}, \dots, \boldsymbol{\pi}) - V_{P,R}^H(\Lambda_P^H(\mu_0, \boldsymbol{\pi}), \boldsymbol{\pi})| \\ &= \left| \mathbb{E}_N \left[\sum_{h=0}^{H-1} \sum_{s \in \mathcal{S}} \hat{\mu}_h(s) \bar{R}(s, \pi_h(s), \hat{\mu}_h) \right] - \sum_{h=0}^{H-1} \sum_{s \in \mathcal{S}} \mu_h(s) R(s, \pi_h(s), \mu_h) \right| \\ &\leq \mathbb{E}_N \left[\sum_{h=0}^{H-1} \left| \sum_{s \in \mathcal{S}} (\hat{\mu}_h(s) \bar{R}(s, \pi_h(s), \hat{\mu}_h) - \mu_h(s) R(s, \pi_h(s), \mu_h)) \right| \right] \\ &\leq \mathbb{E}_N \left[\sum_{h=0}^{H-1} \left(\frac{L_s}{2} \|\mu_h - \hat{\mu}_h\|_1 + L_\mu \|\mu_h - \hat{\mu}_h\|_1 \right) \right]. \end{aligned}$$

The statement of the lemma follows by an application of Lemma A.7. \square

Step 3: Bounding difference in policy deviation. Finally, to conclude the proof of the main theorem of this section, we will prove that the improvement in expectation due to single-sided policy changes are at most of order $\mathcal{O}\left(\frac{1}{\sqrt{N}}\right)$.

Lemma A.9. Suppose $\boldsymbol{\pi} = \{\pi_h\}_{h=0}^{H-1} \in \Pi^H$ and $\boldsymbol{\pi}' = \{\pi'_h\}_{h=0}^{H-1} \in \Pi^H$ arbitrary policies, and $\mu^\boldsymbol{\pi} := \Lambda_P^H(\mu_0, \boldsymbol{\pi})$ is the population distribution induced by $\boldsymbol{\pi}$. Then

$$\begin{aligned} &\left| J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}', \boldsymbol{\pi}, \dots, \boldsymbol{\pi}) - V_{P,R}^H(\Lambda_P^H(\mu_0, \boldsymbol{\pi}), \boldsymbol{\pi}') \right| \\ &\leq \sum_{h=0}^{H-1} \left(\frac{L_\mu}{2} \mathbb{E} [\|\hat{\mu}_h - \mu_h^\boldsymbol{\pi}\|_1] + K_\mu \sum_{h'=0}^{h-1} \mathbb{E} [\|\hat{\mu}_{h'} - \mu_{h'}^\boldsymbol{\pi}\|_1] \right). \end{aligned}$$

825 *Proof.* Define the random variables $\{s_h^i, a_h^i\}_{i,h}, \{\widehat{\mu}_h\}_h$ as in the definition of N -FH-SAG (Definition 3.1). In addition,
 826 define the random variables $\{s_h, a_h\}_h$ evolving according to the FH-MFG with population $\mu^\pi := \{\mu_h^\pi\}_h := \Lambda_P^H(\mu_0, \pi)$
 827 and representative policy π' , independent from the random variables $\{s_h^i, a_h^i\}_{i,h}$. Hence $s_0 \sim \mu_0, a_h \sim \pi'(\cdot|s_h), s_{h+1} \sim$
 828 $P(\cdot|s_h, a_h, \mu_h^\pi)$. Define also for simplicity

$$829 \quad E_N := \left| J_{P,R}^{H,N,(1)}(\pi', \pi, \dots, \pi) - V_{P,R}^H(\Lambda_P^H(\mu_0, \pi), \pi') \right|.$$

830 With these definitions, we have

$$831 \quad E_N = \left| \mathbb{E} \left[\sum_{h=0}^{H-1} R(s_h, a_h, \mu_h^\pi) - \sum_{h=0}^{H-1} R(s_h^1, a_h^1, \widehat{\mu}_h) \right] \right|$$

$$832 \quad \leq \sum_{h=0}^{H-1} \left| \mathbb{E} [R(s_h, a_h, \mu_h^\pi) - R(s_h^1, a_h^1, \widehat{\mu}_h)] \right|. \quad (4)$$

833 Furthermore, for any $h \in \{0, \dots, H-1\}$,

$$834 \quad \left| \mathbb{E} [R(s_h, a_h, \mu_h^\pi) - R(s_h^1, a_h^1, \widehat{\mu}_h)] \right|$$

$$835 \quad \leq \left| \mathbb{E} [R(s_h, a_h, \mu_h^\pi) - R(s_h^1, a_h^1, \mu_h^\pi)] \right|$$

$$836 \quad \quad + \left| \mathbb{E} [R(s_h^1, a_h^1, \mu_h^\pi) - R(s_h^1, a_h^1, \widehat{\mu}_h)] \right|$$

$$837 \quad \leq \left| \mathbb{E} [R(s_h, \pi'_h(s_h), \mu_h^\pi) - R(s_h^1, \pi'_h(s_h^1), \mu_h^\pi)] \right|$$

$$838 \quad \quad + L_\mu \mathbb{E} [\|\mu_h^\pi - \widehat{\mu}_h\|_1]$$

$$839 \quad \leq \frac{1}{2} \|\mathbb{P}[s_h = \cdot] - \mathbb{P}[s_h^1 = \cdot]\|_1 + L_\mu \mathbb{E} [\|\mu_h^\pi - \widehat{\mu}_h\|_1],$$

840 where the last line follows since R is bounded in $[0, 1]$. Replacing this in Equation (4),

$$841 \quad E_N \leq \frac{1}{2} \sum_h \|\mathbb{P}[s_h = \cdot] - \mathbb{P}[s_h^1 = \cdot]\|_1 + L_\mu \sum_h \mathbb{E} [\|\mu_h^\pi - \widehat{\mu}_h\|_1]. \quad (5)$$

842 The first sum above we upper bound in the rest of the proof inductively.

843 Firstly, by definitions of N -FH-SAG and FH-MFG, both s_0^1 and s_0 have distribution μ_0 , hence $\|\mathbb{P}[s_0 = \cdot] - \mathbb{P}[s_0^1 = \cdot]\|_1 = 0$.
 844 Assume that $h \geq 1$. We note that P takes values in Δ_S and the random vector $\widehat{\mu}_h$ takes values in the discrete set
 845 $\{\frac{1}{N}u : u \in \{0, \dots, N\}^S, \sum_s u(s) = N\} \subset \Delta_S$, hence we have the bounds:

$$846 \quad \|\mathbb{P}[s_{h+1} = \cdot] - \mathbb{P}[s_{h+1}^1 = \cdot]\|_1$$

$$847 \quad \leq \left\| \sum_{s,\mu} P(s, \pi'_h(s), \mu) \mathbb{P}[s_h^1 = s, \widehat{\mu}_h = \mu] - \sum_s P(s, \pi'_h(s), \mu_h^\pi) \mathbb{P}[s_h = s] \right\|_1$$

$$848 \quad \leq \left\| \sum_s P(s, \pi'_h(s), \mu_h^\pi) \mathbb{P}[s_h^1 = s] - \sum_s P(s, \pi'_h(s), \mu_h^\pi) \mathbb{P}[s_h = s] \right\|_1$$

$$849 \quad \quad + \left\| \sum_{s,\mu} (P(s, \pi'_h(s), \mu) - P(s, \pi'_h(s), \mu_h^\pi)) \mathbb{P}[s_h^1 = s, \widehat{\mu}_h = \mu] \right\|_1$$

$$850 \quad \leq \|\mathbb{P}[s_h^1 = \cdot] - \mathbb{P}[s_h = \cdot]\|_1 + \sum_{s,\mu} K_\mu \|\mu - \mu_h^\pi\|_1 \mathbb{P}[s_h^1 = s, \widehat{\mu}_h = \mu]$$

$$851 \quad \leq \|\mathbb{P}[s_h^1 = \cdot] - \mathbb{P}[s_h = \cdot]\|_1 + K_\mu \mathbb{E} [\|\widehat{\mu}_h^\pi - \mu_h^\pi\|_1]$$

852 where the last two lines follow from the fact that P is K_μ Lipschitz in μ and stochastic matrices are non-expansive in the
 853 total-variation norm over probability distributions. By induction, we conclude that for all $h \geq 0$, it holds that:

$$854 \quad \|\mathbb{P}[s_h = \cdot] - \mathbb{P}[s_h^1 = \cdot]\|_1 \leq K_\mu \sum_{h'=0}^h \mathbb{E} [\|\widehat{\mu}_{h'}^\pi - \mu_{h'}^\pi\|_1].$$

Placing this result into Equation (5), we obtain the statement of the lemma. \square

Since $\mathbb{E}[\|\widehat{\mu}_{h'} - \mu_{h'}^{\pi'}\|_1]$ above in the theorem is of the order of $\mathcal{O}(1/\sqrt{N})$ by the result in step 1, the result above allows us to bound exploitability in the N -FH-SAG.

Conclusion and Statement of Result. Finally, we can merge the results up until this stage to upper bound the exploitability. By definition of the FH-MFG-NE, we have:

$$\delta \geq \max_{\pi' \in \Pi^H} V_{P,R}^H(\Lambda_P^H(\mu_0, \pi_\delta), \pi') - V_{P,R}^H(\Lambda_P^H(\mu_0, \pi_\delta), \pi_\delta)$$

The upper bounds on the deviation between $V_{P,R}^H$ and $J_{P,R}^{H,N,(1)}$ from the previous steps directly yields the statement of the theorem. We state it below for completeness.

Theorem A.10. *It holds that*

$$\mathcal{E}_{P,R}^{H,N,(1)}(\pi_\delta, \dots, \pi_\delta) \leq 2\delta + \frac{C_1}{\sqrt{N}} + \frac{C_2}{N} = O\left(\delta + \frac{1}{\sqrt{N}}\right)$$

where π_δ is a δ -FH-MFG Nash equilibrium and

$$C_1 = |\mathcal{S}| \sqrt{\frac{\pi}{2}} \left((2L_\mu + \frac{L_s}{2}) \sum_{h=0}^{H-1} \frac{1 - L_{pop,\mu}^{h+1}}{1 - L_{pop,\mu}} + K_\mu \sum_{h=0}^{H-1} \sum_{i=0}^{h-1} \frac{1 - L_{pop,\mu}^{i+1}}{1 - L_{pop,\mu}} \right)$$

$$C_2 = L_\mu K_a \sum_{h=0}^{H-1} \frac{1 - L_{pop,\mu}^h}{1 - L_{pop,\mu}} + K_a K_\mu \sum_{h=0}^{H-1} \sum_{i=0}^{h-1} \frac{1 - L_{pop,\mu}^i}{1 - L_{pop,\mu}},$$

where we use shorthand notation $\frac{1 - L_{pop,\mu}^k}{1 - L_{pop,\mu}} := k - 1$ when $L_{pop,\mu} = 1$.

A note on constants. Note that constants C_1, C_2 in Theorem A.10 depend on horizon with $\frac{H^2}{1 - L_{pop,\mu}}$ if $L_{pop,\mu} < 1$, with H^3 if $L_{pop,\mu} = 1$ and with $H^2 \frac{1 - L_{pop,\mu}^{H+1}}{1 - L_{pop,\mu}}$ if $L_{pop,\mu} > 1$.

A.3. Lower Bound for FH-MFG: Extended Proof of Theorem 3.3

The proof will be by construction: we will explicitly define an FH-MFG where the optimal policy for the N -agent game diverges quickly from the FH-MFG-NE policy.

Preliminaries. We first define a few utility functions. Define $\mathbf{g} : \Delta_2 \rightarrow B_{\infty,+}^2 := \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_\infty = 1, x_1, x_2 \geq 0\}$ and $\mathbf{h} : \Delta_2 \rightarrow [0, 1]^2$ as follows:

$$\mathbf{g}(x_1, x_2) := \begin{pmatrix} \mathbf{g}_1(x_1, x_2) \\ \mathbf{g}_2(x_1, x_2) \end{pmatrix} := \begin{pmatrix} \frac{x_1}{\max\{x_1, x_2\}} \\ \frac{x_2}{\max\{x_1, x_2\}} \end{pmatrix},$$

$$\mathbf{h}(x_1, x_2) := \begin{pmatrix} \mathbf{h}_1(x_1, x_2) \\ \mathbf{h}_2(x_1, x_2) \end{pmatrix} := \begin{pmatrix} \max\{4x_2, 1\} \\ \max\{4x_1, 1\} \end{pmatrix}.$$

Furthermore, for any $\epsilon > 0$ we define $\omega_\epsilon : [0, 1] \rightarrow [0, 1]$ as:

$$\omega_\epsilon(x) = \begin{cases} 1, & x > 1/2 + \epsilon \\ 0, & x < 1/2 - \epsilon \\ \frac{1}{2} + \frac{x - 1/2}{2\epsilon}, & x \in [1/2 - \epsilon, 1/2 + \epsilon] \end{cases}.$$

$\epsilon \in (0, 1/2)$ will be specified later.

It is straightforward to verify that \mathbf{g} has an inverse in its domain given by

$$\mathbf{g}^{-1}(x_1, x_2) = \left(\frac{x_1}{x_1 + x_2}, \frac{x_2}{x_1 + x_2} \right), \forall (x_1, x_2) \in B_{\infty,+}^2.$$

Furthermore, it holds for $\mathbf{x} = (x_1, x_2) \in B_{\infty,+}^2, \mathbf{y} = (y_1, y_2) \in B_{\infty,+}^2$

$$\begin{aligned} & \|\mathbf{g}^{-1}(\mathbf{x}) - \mathbf{g}^{-1}(\mathbf{y})\|_1 \\ &= \left| \frac{x_1}{x_1 + x_2} - \frac{y_1}{y_1 + y_2} \right| + \left| \frac{x_2}{x_1 + x_2} - \frac{y_2}{y_1 + y_2} \right| \\ &= \left| \frac{x_1(y_2 - x_2) + x_2(x_1 - y_1)}{(x_1 + x_2)(y_1 + y_2)} \right| + \left| \frac{x_2(y_1 - x_1) + x_1(x_2 - y_2)}{(x_1 + x_2)(y_1 + y_2)} \right| \\ &\leq 2\|\mathbf{x} - \mathbf{y}\|_1, \end{aligned}$$

and likewise for $\mathbf{u}, \mathbf{v} \in \Delta_2$, letting $u_+ := \max\{u_1, u_2\}, v_+ := \max\{v_1, v_2\}$,

$$\begin{aligned} \|\mathbf{g}(\mathbf{u}) - \mathbf{g}(\mathbf{v})\|_1 &= \left| \frac{u_1}{u_+} - \frac{v_1}{v_+} \right| + \left| \frac{u_2}{u_+} - \frac{v_2}{v_+} \right| \\ &= \left| \frac{u_1 v_+ - v_1 u_+}{u_+ v_+} \right| + \left| \frac{u_2 v_+ - u_+ v_2}{u_+ v_+} \right| \leq 2\|\mathbf{u} - \mathbf{v}\|_1. \end{aligned}$$

This follows from considering cases and observation that $u_+ \geq 1/2, v_+ \geq 1/2$. Then for all $\mathbf{u}, \mathbf{v} \in \Delta_2$, \mathbf{g}, \mathbf{h} have the bi-Lipschitz and Lipschitz properties:

$$\frac{1}{2}\|\mathbf{u} - \mathbf{v}\|_1 \leq \|\mathbf{g}(\mathbf{u}) - \mathbf{g}(\mathbf{v})\|_1 \leq 2\|\mathbf{u} - \mathbf{v}\|_1, \quad (6)$$

$$\|\mathbf{h}(\mathbf{u}) - \mathbf{h}(\mathbf{v})\|_1 \leq 4\|\mathbf{u} - \mathbf{v}\|_1. \quad (7)$$

Likewise, ω_ϵ , being piecewise linear, also satisfies the Lipschitz condition: $|\omega_\epsilon(x) - \omega_\epsilon(y)| \leq \frac{1}{2\epsilon}|x - y|, \forall x, y \in [0, 1]$.

Defining the FH-MFG. We take a particular FH-MFG with 6 states, 2 actions. Define the state-actions sets:

$$\mathcal{S} = \{s_{\text{Left}}, s_{\text{Right}}, s_{\text{LA}}, s_{\text{LB}}, s_{\text{RA}}, s_{\text{RB}}\}, \quad \mathcal{A} = \{a_A, a_B\}.$$

Intuitively, the ‘‘main’’ states of the game are $s_{\text{Left}}, s_{\text{Right}}$ and the 4 states $s_{\text{LA}}, s_{\text{LB}}, s_{\text{RA}}, s_{\text{RB}}$ are dummy states that keep track of which actions were taken by which percentage of players used to introduce a dependency of the rewards on the distribution of agents over actions as well as states. Define the initial probabilities μ_0 by:

$$\begin{aligned} \mu_0(s_{\text{Left}}) &= \mu_0(s_{\text{Right}}) = 1/2, \\ \mu_0(s_{\text{LA}}) &= \mu_0(s_{\text{RA}}) = \mu_0(s_{\text{LB}}) = \mu_0(s_{\text{RB}}) = 0. \end{aligned}$$

When at the states $s_{\text{Left}}, s_{\text{Right}}$, the transition probabilities are defined for all $\mu \in \Delta_{\mathcal{S}}$ by:

$$\begin{aligned} P(s_{\text{LA}}|s_{\text{Left}}, a_A, \mu) &= 1, & P(s_{\text{LB}}|s_{\text{Left}}, a_B, \mu) &= 1, \\ P(s_{\text{RA}}|s_{\text{Right}}, a_A, \mu) &= 1, & P(s_{\text{RB}}|s_{\text{Right}}, a_B, \mu) &= 1. \end{aligned}$$

That is, the agent transitions to one of $\{s_{\text{LA}}, s_{\text{RA}}, s_{\text{RB}}, s_{\text{LB}}\}$ to remember its last action and left-right state. When at states $\{s_{\text{LA}}, s_{\text{RA}}, s_{\text{RB}}, s_{\text{LB}}\}$, the transition probabilities are:

$$\begin{aligned} & \text{If } s \in \{s_{\text{LA}}, s_{\text{LB}}, s_{\text{RA}}, s_{\text{RB}}\} : \\ & P(s'|s, a, \mu) = \begin{cases} \omega_\epsilon(\mu(s_{\text{LA}}) + \mu(s_{\text{LB}})), & \text{if } s' = s_{\text{Left}} \\ \omega_\epsilon(\mu(s_{\text{RA}}) + \mu(s_{\text{RB}})), & \text{if } s' = s_{\text{Right}} \end{cases}, \forall \mu, a. \end{aligned}$$

The other non-defined transition probabilities are of course 0.

Finally, let $\alpha, \beta > 0$ such that $\alpha + \beta < 1$ (to be also defined later). The reward functions are defined for all $\mu \in \Delta_{\mathcal{S}}$ as follows:

$$\begin{aligned}
 R(s_{\text{Left}}, a_A, \mu) &= R(s_{\text{Left}}, a_B, \mu) = 0, \\
 R(s_{\text{Right}}, a_A, \mu) &= R(s_{\text{Right}}, a_B, \mu) = 0, \\
 \begin{pmatrix} R(s_{\text{LA}}, a_A, \mu) \\ R(s_{\text{LB}}, a_A, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(\mu(s_{\text{LA}}) + \mu(s_{\text{LB}}), \mu(s_{\text{RA}}) + \mu(s_{\text{RB}})) \\
 &\quad + \alpha \mathbf{h}(\mu(s_{\text{LA}}), \mu(s_{\text{LB}})) \\
 \begin{pmatrix} R(s_{\text{LA}}, a_B, \mu) \\ R(s_{\text{LB}}, a_B, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(\mu(s_{\text{LA}}) + \mu(s_{\text{LB}}), \mu(s_{\text{RA}}) + \mu(s_{\text{RB}})) \\
 &\quad + \alpha \mathbf{h}(\mu(s_{\text{LA}}), \mu(s_{\text{LB}})) + \beta \mathbf{1} \\
 \begin{pmatrix} R(s_{\text{RA}}, a_A, \mu) \\ R(s_{\text{RB}}, a_A, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(\mu(s_{\text{RA}}) + \mu(s_{\text{RB}}), \mu(s_{\text{LA}}) + \mu(s_{\text{LB}})) \\
 &\quad + \alpha \mathbf{h}(\mu(s_{\text{RA}}), \mu(s_{\text{RB}})) \\
 \begin{pmatrix} R(s_{\text{RA}}, a_B, \mu) \\ R(s_{\text{RB}}, a_B, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(\mu(s_{\text{RA}}) + \mu(s_{\text{RB}}), \mu(s_{\text{LA}}) + \mu(s_{\text{LB}})) \\
 &\quad + \alpha \mathbf{h}(\mu(s_{\text{RA}}), \mu(s_{\text{RB}})) + \beta \mathbf{1}
 \end{aligned}$$

Note that only at odd steps do the agents get a reward, and at this step, it does not matter which action the agent plays, only the state among $\{s_{\text{LA}}, s_{\text{LB}}, s_{\text{RA}}, s_{\text{RB}}\}$ and the population distribution. The parameters ϵ, α, β of the above FH-MFG are “free” parameters to be specified later. We visualize the FH-MFG in Figure 1.

A minor remark. The arguments of \mathbf{g} above will be with probability one in the set Δ_2 at odd-numbered time steps, but to formally satisfy the Lipschitz condition $R \in \mathcal{R}_2$ one can for instance replace $\mathbf{g}(\mu(s_{\text{RA}}) + \mu(s_{\text{RB}}), \mu(s_{\text{LA}}) + \mu(s_{\text{LB}}))$ with $\mathbf{g}(\mu(s_{\text{RA}}) + \mu(s_{\text{RB}}) + \mu(s_{\text{Left}}), \mu(s_{\text{LA}}) + \mu(s_{\text{LB}}) + \mu(s_{\text{Right}}))$ in the definitions, which will not impact the analysis since at odd timesteps $\mu(s_{\text{Right}}) = \mu(s_{\text{Left}}) = 0$ for both the FH-MFG and N -FH-SAG.

Note that with these definitions, $P \in \mathcal{P}_{1/2\epsilon}, R \in \mathcal{R}_2$ since only $\forall s, s' \in \mathcal{S}, a, a' \in \mathcal{A}, \mu, \mu' \in \Delta_{\mathcal{S}}$, we have by the definitions:

$$\|P(\cdot | s, a, \mu) - P(\cdot | s', a', \mu')\|_1 \leq 2d(s, s') + 2d(a, a') + \frac{1}{2\epsilon} \|\mu - \mu'\|_1, \quad (8)$$

$$|R(s, a, \mu) - R(s', a', \mu')| \leq d(s, s') + d(a, a') + 2\|\mu - \mu'\|_1, \quad (9)$$

for any $\alpha, \beta > 0$ with $\alpha + \beta < 1$ and $\alpha < \frac{1}{4}$, using the Lipschitz conditions in (6), (7).

Step 1: Solution of the FH-MFG. Next, we solve the infinite player FH-MFG and show that the policy $\pi_H^* := \{\pi_h^*\}_{h=0}^{H-1}$ given by:

$$\pi_h^*(a | s) := \begin{cases} 1, & \text{if } h \text{ odd and } a = a_B \\ \frac{1}{2}, & \text{if } h \text{ even} \\ 0, & \text{if } h \text{ odd and } a = a_A \end{cases}$$

It is easy to verify in this case that, if $\mu^* := \{\mu_h^*\}_h$ is induced by π^* :

$$\begin{aligned}
 \mu_h^*(s_{\text{LA}}) &= \mu_h^*(s_{\text{LB}}) = \mu_h^*(s_{\text{RA}}) = \mu_h^*(s_{\text{RB}}) = 1/4, & \text{if } h \text{ odd,} \\
 \mu_h^*(s_{\text{Left}}) &= \mu_h^*(s_{\text{Right}}) = 1/2, & \text{if } h \text{ even.}
 \end{aligned}$$

In this case, the induced rewards in odd steps are state-independent (it is the same for all states $s_{\text{RA}}, s_{\text{RB}}, s_{\text{LA}}, s_{\text{LB}}$), therefore the policy π^* is the optimal best response to the population and a FH-MFG.

In fact, π^* is unique up to modifications in zero-probability sets (e.g., modifying $\pi_h^*(s_{\text{Left}})$ for odd h , for which $\mathbb{P}[s_h = s_{\text{Left}}] = 0$). To see this, for any policy $\pi \in \Pi_H$, it holds that

$$\begin{aligned}
 \mu_h^\pi(s_{\text{Left}}) &= \mu_h^\pi(s_{\text{Right}}) = 1/2, & \text{if } h \text{ even,} \\
 \mu_h^\pi(s_{\text{LA}}) + \mu_h^\pi(s_{\text{LB}}) &= \mu_h^\pi(s_{\text{RA}}) + \mu_h^\pi(s_{\text{RB}}) = 1/2, & \text{if } h \text{ odd,}
 \end{aligned}$$

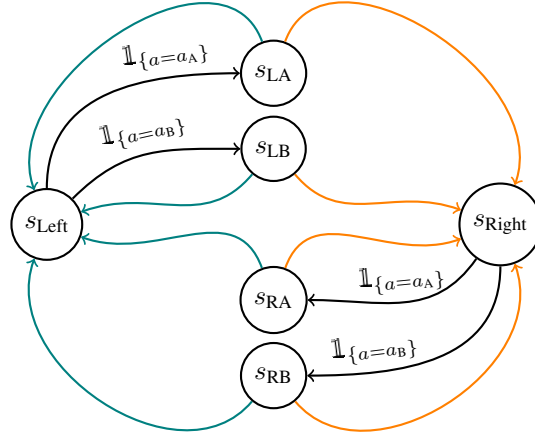


Figure 1. Visualization of the counterexample. All orange edges have probability $\omega_\varepsilon(\mu(s_{RA}) + \mu(s_{RB}))$, green edges have probability $\omega_\varepsilon(\mu(s_{LA}) + \mu(s_{LB}))$ independent of action taken. Edges with probability 0 are not drawn.

as the action of the agent does not affect transition probabilities between $s_{\text{Left}}, s_{\text{Right}}$ in even rounds. Moreover, as odd stages, the action rewards terms only depend on the state apart from the positive additional term $\beta \mathbf{1}$, so the only optimal action will be a_B . Finally, for $\alpha > 0$, the actions a_A, a_B must be played with equal probability as otherwise the term $\alpha \mathbf{h}(\mu(s_{RA}), \mu(s_{RB}))$ will lead to the action with lower probability assigned by being optimal.

Step 2: Population divergence in N -FH-MFG. We will analyze the empirical population distribution deviation from μ^* , namely, we will lower bound $\mathbb{E}[\|\mu_h^* - \hat{\mu}_h\|_1]$. The results in this step will be valid for *any* policy profile $(\pi^1, \dots, \pi^N) \in \Pi$: we emphasize that at even h , $\hat{\mu}_h$ is independent of agent policies in the N player game. In this step, we also fix $1/2\varepsilon = 8$.

We will analyze $\hat{\mu}_h$ at all even steps $h = 2m$ where $m \in \mathbb{N}_{\geq 0}$. Define the sequence of random variables for all $m \in \mathbb{N}_{\geq 0}$ as $X_m := \hat{\mu}_{2m}(s_{\text{Left}})$. Define $\mathcal{G} := \{\frac{k}{N} : k = 0, \dots, N\}$. Note that for all even $h = 2m$, it holds almost surely that $\hat{\mu}_h(s_{\text{Left}}), \hat{\mu}_h(s_{\text{Right}}) \in \mathcal{G}$. By the definition of the MFG, it holds for any $m \geq 0, k \in [N]$ that

$$\begin{aligned} \mathbb{P}[NX_0 = k] &= \binom{N}{k} 2^{-N}, \\ \mathbb{P}[NX_{m+1} = k | X_m] &= \binom{N}{k} (\omega_\varepsilon(X_m))^k (1 - \omega_\varepsilon(X_m))^{N-k}, \end{aligned}$$

that is, given X_m , NX_{m+1} is binomially distributed with $NX_{m+1} \sim \text{Binom}(N, \omega_\varepsilon(X_m))$ without any dependence on the actions played by agents. Therefore

$$\mathbb{E}[X_{m+1} | X_m] = \omega_\varepsilon(X_m), \quad \text{Var}[X_{m+1} | X_m] \leq \frac{1}{4N}.$$

We define the following set $\mathcal{G}_* := \{0, 1\} \subset \mathcal{G}$. By the definition of the mechanics, if $x \in \mathcal{G}_*, m \in \mathbb{N}_{\geq 0}$, it holds for all $m' > m$ that $\mathbb{P}[X_{m'} = x | X_m = x] = 1$, that is once the Markovian random process X_m hits \mathcal{G}_* , it will remain in \mathcal{G}_* . Furthermore, for $K := \lfloor \log_5 \sqrt{N} \rfloor$, and for $k = 0, \dots, K$ define the level sets:

$$\mathcal{G}_{-1} := \mathcal{G}, \quad \mathcal{G}_k := \left\{ x \in \mathcal{G} : \left| x - \frac{1}{2} \right| \geq \frac{5^k}{2\sqrt{N}} \right\}.$$

For all $k \geq K$, define $\mathcal{G}_k := \mathcal{G}_*$.

1100 Firstly, we have that

$$\begin{aligned}
 1101 \quad \mathbb{P}[X_0 \in \mathcal{G}_0] &= \mathbb{P} \left[\left| \frac{1}{N} \sum_i \mathbb{1}_{\{s_0^i = s_{\text{Left}}\}} - \frac{1}{2} \right| \geq \frac{1}{2\sqrt{N}} \right] \\
 1102 &= \mathbb{P} \left[\left| \sum_i \mathbb{1}_{\{s_0^i = s_{\text{Left}}\}} - \frac{N}{2} \right| \geq \frac{\sqrt{N}}{2} \right] \geq \frac{1}{10}, \\
 1103 & \\
 1104 & \\
 1105 & \\
 1106 &
 \end{aligned}$$

1107 where in the last line we applied the anti-concentration result of Lemma A.4 on the sum of independent Bernoulli random
 1108 variables $\mathbb{1}_{\{s_0^i = s_{\text{Left}}\}}$ for $i \in [N]$.

1110 Next, assume that for some $m \in 1, \dots, K-1$ we have $p \in \mathcal{G}_m$. If $\omega_\epsilon(p) \in \{0, 1\}$, it holds trivially that $\mathbb{P}[X_{m+1} \in$
 1111 $\mathcal{G}_{m+1} | X_m = p] = 1$. Otherwise, if $\omega_\epsilon(p) \in (0, 1)$,

$$\begin{aligned}
 1112 \quad &\mathbb{P}[X_{m+1} \in \mathcal{G}_{m+1} | X_m = p] \\
 1113 &= \mathbb{P} \left[\left| X_{m+1} - \frac{1}{2} \right| \geq \frac{5^{m+1}}{2\sqrt{N}} \mid X_m = p \right] \\
 1114 &\geq \mathbb{P} \left[\left| \omega_\epsilon(p) - \frac{1}{2} \right| - |X_{m+1} - \omega_\epsilon(p)| \geq \frac{5^{m+1}}{2\sqrt{N}} \mid X_m = p \right]. \\
 1115 & \\
 1116 & \\
 1117 & \\
 1118 &
 \end{aligned}$$

1119 Since in this case $|\omega_\epsilon(X_m) - \frac{1}{2}| = |\omega_\epsilon(X_m) - \omega_\epsilon(\frac{1}{2})| \geq 1/2\epsilon |X_m - \omega_\epsilon(\frac{1}{2})|$, we have

$$\begin{aligned}
 1120 \quad &\mathbb{P}[X_{m+1} \in \mathcal{G}_{m+1} | X_m = p] \\
 1121 &\geq \mathbb{P} \left[\left| \omega_\epsilon(p) - \frac{1}{2} \right| - |X_{m+1} - \omega_\epsilon(p)| \geq \frac{5^{m+1}}{2\sqrt{N}} \mid X_m = p \right] \\
 1122 &= \mathbb{P} \left[|X_{m+1} - \omega_\epsilon(p)| \leq \left| \omega_\epsilon(p) - \frac{1}{2} \right| - \frac{5^{m+1}}{2\sqrt{N}} \mid X_m = p \right] \\
 1123 &\geq \mathbb{P} \left[|X_{m+1} - \omega_\epsilon(p)| \leq 8 \frac{5^m}{2\sqrt{N}} - \frac{5^{m+1}}{2\sqrt{N}} \mid X_m = p \right] \\
 1124 &= \mathbb{P} \left[|X_{m+1} - \omega_\epsilon(p)| \leq 3 \frac{5^m}{2\sqrt{N}} \mid X_m = p \right] \\
 1125 &\geq 1 - 2 \exp \left\{ -\frac{9}{50} 25^{m+1} \right\} \\
 1126 & \\
 1127 & \\
 1128 & \\
 1129 & \\
 1130 & \\
 1131 & \\
 1132 & \\
 1133 &
 \end{aligned}$$

1134 where in the last line we invoked the Hoeffding concentration bound (Lemma A.2).

1135 Using the above result inductively for $m \in 0, \dots, K$ it holds that

$$\begin{aligned}
 1136 \quad \mathbb{P}[X_m \in \mathcal{G}_m | X_0 \in \mathcal{G}_0] &\geq \prod_{m'=1}^m \mathbb{P}[X_{m'} \in \mathcal{G}_{m'} | X_{m'-1} \in \mathcal{G}_{m'-1}] \\
 1137 &\geq \prod_{m'=1}^m \left(1 - 2 \exp \left\{ -\frac{9}{50} 25^{m'} \right\} \right) \\
 1138 &\geq \left(1 - 2 \sum_{m'=0}^{\infty} \exp \left\{ -\frac{9}{50} 25^{m'+1} \right\} \right) \\
 1139 &\geq \left(1 - 2 \sum_{m'=0}^{\infty} \exp \left\{ -\frac{9}{2} m' - \frac{9}{2} \right\} \right) \\
 1140 &\geq \left(1 - \frac{2e^{-9/2}}{1 - e^{-9/2}} \right) \geq \frac{9}{10}. \\
 1141 & \\
 1142 & \\
 1143 & \\
 1144 & \\
 1145 & \\
 1146 & \\
 1147 & \\
 1148 & \\
 1149 & \\
 1150 &
 \end{aligned}$$

1151 Since for $k > K$, $\mathbb{P}[X_{k+1} \in \mathcal{G}_* | X_k \in \mathcal{G}_*] = 1$ and $\mathbb{P}[X_0 \in \mathcal{G}_0] \geq 1/10$, it also holds that

$$1152 \quad \mathbb{P}[X_m \in \mathcal{G}_m, \forall m \geq 0] \geq \frac{9}{100}.$$

1154

1155 Finally, we use the above lower bound on the probability to lower bound the expectation:

$$\begin{aligned}
 1156 & \\
 1157 & \mathbb{E} [\|\widehat{\mu}_{2m} - \mu_{2m}\|_1] \geq \mathbb{P}[X_m \in \mathcal{G}_m] \mathbb{E} [\|\widehat{\mu}_{2m} - \mu_{2m}\|_1 | X_m \in \mathcal{G}_m] \\
 1158 & \geq \mathbb{P}[X_m \in \mathcal{G}_m] \mathbb{E} [2|X_m - 1/2| | X_m \in \mathcal{G}_m] \\
 1159 & \geq \frac{9}{100} \min \left\{ \frac{5^m}{\sqrt{N}}, 1 \right\}. \\
 1160 & \\
 1161 &
 \end{aligned}$$

1162 For odd $h = 2m + 1$, we also have the inequality

$$\begin{aligned}
 1163 & \\
 1164 & \mathbb{E} [\|\widehat{\mu}_{2m+1} - \mu_{2m+1}\|_1] \geq \mathbb{E} [\|\widehat{\mu}_{2m} - \mu_{2m}\|_1] \\
 1165 & \geq \frac{9}{100} \min \left\{ \frac{5^m}{\sqrt{N}}, 1 \right\}. \\
 1166 & \\
 1167 &
 \end{aligned}$$

1169 which completes the first statement of the theorem (as $5^{H/2} = \Omega(2^H)$).

1170 **Step 3: Hitting time for \mathcal{G}_* .** We will show that the empirical distribution of agent states almost always concentrates on one of $s_{\text{Left}}, s_{\text{Right}}$ during the even rounds in the N -player game, and bound the expected waiting time for this to happen. The distributions of agents over states $s_{\text{Left}}, s_{\text{Right}}$ in the even rounds are policy independent (they are not affected by which actions are played): hence the results from Step 2 still hold for the population distribution and the expected time computed in this step will be valid for any policy.

1176 For simplicity, we define the FH-MFG for the non-terminating infinite horizon chain, and we will compute value functions up to horizon H . Define the (random) hitting time τ as follows:

$$1177 \tau := \inf \{m \geq 0 : \widehat{\mu}_{2m}(s_{\text{Left}}) \in \mathcal{G}_*\} = \inf \{m \geq 0 : X_m \in \mathcal{G}_*\}.$$

1181 Note that for any $p \in \mathcal{G}$, it holds that $\mathbb{P}[X_{m+1} \in \mathcal{G}_* | X_m = p] = \widehat{\mu}_{2m}(s_{\text{Left}})^N + \widehat{\mu}_{2m}(s_{\text{Right}})^N = p^N + (1-p)^N \geq 2^{-N}$. Therefore for all m it holds that $\mathbb{P}[\widehat{\mu}_{2m} \notin \mathcal{G}_*] \leq (1 - 2^{-N})^{m-1}$. By the Borel-Cantelli lemma, we can conclude that $\tau < \infty$ almost surely, and in particular $T_\tau := \mathbb{E}[\tau | X_0 = x] < \infty$ for any $x \in \mathcal{G}$.

1185 Next, we compute the expected value T_τ . Define the following two quantities:

$$\begin{aligned}
 1186 & \\
 1187 & T_{-1} := \sup_{x \in \mathcal{G}_{-1}} \{\mathbb{E}[\tau | X_0 = x]\} \\
 1188 & \\
 1189 & T_0 := \sup_{x \in \mathcal{G}_0} \{\mathbb{E}[\tau | X_0 = x]\}. \\
 1190 & \\
 1191 &
 \end{aligned}$$

1192 First, we compute an upper bound for T_0 . Define the event:

$$1193 E_0 := \bigcap_{m' \in [K]} \{X_{m'} \in \mathcal{G}_{m'}\}.$$

1197 Then, T_0 is upper bounded by:

$$\begin{aligned}
 1198 & \\
 1199 & T_0 = \sup_{x \in \mathcal{G}_0} \mathbb{E}[\tau | X_0 = x] \\
 1200 & = \sup_{x \in \mathcal{G}_0} \mathbb{E}[\tau | E_0, X_0 = x] \mathbb{P}[E_0 | X_0 = x] \\
 1201 & \quad + \mathbb{E}[\tau | E_0^c, X_0 = x] \mathbb{P}[E_0^c | X_0 = x] \\
 1202 & \leq \sup_{x \in \mathcal{G}_0} \mathbb{E}[\tau | E_0, X_0 = x] \mathbb{P}[E_0 | X_0 = x] \\
 1203 & \quad + \mathbb{E}[\tau | E_0^c, X_0 = x] \mathbb{P}[E_0^c | X_0 = x] \\
 1204 & \leq K \frac{9}{10} + (K + T_{-1}) \frac{1}{10} = K + \frac{T_{-1}}{10} \\
 1205 & \\
 1206 & \\
 1207 & \\
 1208 & \\
 1209 &
 \end{aligned}$$

1210 where in the last step we used the lower bound on $\mathbb{P}[E_0]$ from Step 2. Similarly for T_{-1} , from the one-sided anti-concentration
 1211 bound (Lemma A.4) it holds that:

$$\begin{aligned} 1212 \quad T_{-1} &\leq \sup_{x \in \mathcal{G}_{-1}} \mathbb{E}[\tau | X_0 = x] \\ 1213 &\leq \mathbb{E}[\tau | x \in \mathcal{G}_0, X_0 = x] \mathbb{P}[x \in \mathcal{G}_0 | X_0 = x] \\ 1214 &\quad + \mathbb{E}[\tau | x \notin \mathcal{G}_0, X_0 = x] \mathbb{P}[x \notin \mathcal{G}_0 | X_0 = x] \\ 1215 &\leq \frac{1}{20}(T_0 + 1) + \frac{19}{20}(T_{-1} + 1), \end{aligned}$$

1216 the last line following since $T_{-1} > T_0$ by definition. Solving the two inequalities, we obtain

$$1217 \quad T_\tau \leq T_{-1} \leq \frac{200}{9} + \frac{10K}{9} \leq 23 + \frac{5}{9} \log_5 N.$$

1218 **Step 4: Ergodic optimal response to N -players.** Next, we formulate a policy $\pi^{\text{br}} = \{\pi_h^{\text{br}}\}_{h=0}^{H-1} \in \Pi^H$ that is ergodically
 1219 optimal for the N -player game and can exploit a population that deploys the unique FH-MFG-NE. For all h , the optimal
 1220 policy will be defined by:

$$1221 \quad \pi_h^{\text{br}}(a|s) = \begin{cases} 1, & \text{if } s = s_{\text{Left}}, a = a_A \\ 1, & \text{if } s = s_{\text{Right}}, a = a_B \\ 1, & \text{if } s \notin \{s_{\text{Left}}, s_{\text{Right}}\}, a = a_B \\ 0, & \text{otherwise} \end{cases}$$

1222 Intuitively, π_h^{br} becomes optimal once all the agents are concentrated in the same states during the even rounds, which
 1223 happens very quickly as shown in Step 3. Assume that agents $i = 2, \dots, N$ deploy the unique FH-MFG-NE $\pi^i = \pi^*$, and for
 1224 agent $i = 1$, $\pi^1 = \pi^{\text{br}}$. We decompose the three components of the rewards for the first agent, as defined in the construction
 1225 of the MFG (Step 1):

$$\begin{aligned} 1226 \quad J_{P,R}^{H,N,(1)}(\pi^{\text{br}}, \pi^*, \dots, \pi^*) &= \mathbb{E} \left[\sum_{\substack{h \text{ odd} \\ 0 \leq h \leq H}} (1 - \alpha - \beta) R_h^{1,\mathbf{g}} + \alpha R_h^{1,\mathbf{h}} + \beta \mathbb{1}_{\{a_h^1 = a_B\}} \right] \\ 1227 &\geq (1 - \alpha - \beta) \mathbb{E} \left[\sum_{\text{odd } h=0}^{H-1} R_h^{1,\mathbf{g}} \right] + \beta \left[\frac{H}{2} \right] \end{aligned}$$

1228 as by definition clearly $\mathbb{E} \left[\mathbb{1}_{\{a_h^1 = a_B\}} \right] = 1$ for all odd h and $R_h^{\mathbf{h}} \geq 0$ almost surely.

1229 We analyze the terms $R_h^{1,\mathbf{g}}$ when the first agent follows π^{br} . By the definition of the dynamics and π^{br} , it holds that

$$1230 \quad R_h^{1,\mathbf{g}} = g_1(\widehat{\mu}_{h-1}(s_{h-1}^1), \widehat{\mu}_{h-1}(\bar{s}_{h-1}^1))$$

1231 where $\bar{s}_{h-1}^1 := s_{\text{Left}}$ if $s_{h-1}^1 = s_{\text{Right}}$ and $\bar{s}_{h-1}^1 := s_{\text{Right}}$ if $s_{h-1}^1 = s_{\text{Left}}$. As $\mathbb{P}[s_{h-1}^1 = \cdot, \dots, s_{h-1}^N = \cdot]$ at even step $h-1$
 1232 is permutation invariant, it holds that $\mathbb{P}[s_{h-1}^1 = \cdot | \widehat{\mu}_{h-1} = \mu] = \mu(\cdot)$ for any $\mu \in \mathcal{G}$. Therefore,

$$\begin{aligned} 1233 \quad \mathbb{E}[R_h^{1,\mathbf{g}}] &= \sum_{\substack{\mu \in \mathcal{G} \\ s \in \{s_{\text{Left}}, s_{\text{Right}}\}}} \mathbb{P}[\widehat{\mu}_{h-1} = \mu] \mathbb{P}[s_{h-1}^1 = s | \widehat{\mu}_{h-1} = \mu] \\ 1234 &\quad \mathbb{E}[R_h^{1,\mathbf{g}} | s_{h-1}^1 = s, \widehat{\mu}_{h-1} = \mu] \\ 1235 &= \sum_{\substack{\mu \in \mathcal{G} \\ s \in \{s_{\text{Left}}, s_{\text{Right}}\}}} \mathbb{P}[\widehat{\mu}_{h-1} = \mu] \mu(s) g_1(\mu(s), \mu(\bar{s})) \geq 1/2, \end{aligned}$$

1265 as for any μ , if s is such that $\mu(s) \geq \mu(\bar{s})$ then $g_1(\mu(s), \mu(\bar{s})) = 1$. Furthermore, by the definition of the hitting time τ ,
 1266 for any odd $h \geq 1$, $\mathbb{E}[R_h^g | 2\tau < h] = \mathbb{E}[R_h^g | \hat{\mu}_{h-1}(s_{\text{Left}}) \in \mathcal{G}_*] = 1$, as after time 2τ the action a_A will be optimal with
 1267 reward $R_h^g = 1$ almost surely, as π^{br} chooses action a_A at even steps.

1268 Finally, using the lower bound of $1/2$ for R_h^g when $h < 2\tau$ and that $R_h^g = 1$ when $h > 2\tau$, we obtain:
 1269

$$\begin{aligned}
 1270 \quad \mathbb{E} \left[\sum_{\substack{h \text{ odd} \\ 0 \leq h \leq H}} R_h^g \right] &= \mathbb{E} \left[\sum_{0 \leq h \leq \min\{2\tau, H\}} R_h^{1,g} + \sum_{\substack{h \text{ odd} \\ \min\{2\tau, H\} + 1 \leq h < H}} R_h^{1,g} \right] \\
 1271 &\geq \mathbb{E} \left[\frac{1}{2} \min \left\{ \tau, \left\lfloor \frac{H}{2} \right\rfloor \right\} + \left(\left\lfloor \frac{H}{2} \right\rfloor - \min \left\{ \tau, \left\lfloor \frac{H}{2} \right\rfloor \right\} \right) \right] \\
 1272 &\geq \left\lfloor \frac{H}{2} \right\rfloor - \frac{1}{2} \mathbb{E} \left[\min \left\{ \tau, \left\lfloor \frac{H}{2} \right\rfloor \right\} \right] \\
 1273 &\geq \left\lfloor \frac{H}{2} \right\rfloor - \frac{1}{2} \mathbb{E}[\tau] = \left\lfloor \frac{H}{2} \right\rfloor - \frac{T_\tau}{2}
 \end{aligned}$$

1274 Merging the inequalities above, we obtain

$$1275 \quad J_{P,R}^{H,N,(1)}(\pi^{\text{br}}, \pi^*, \dots, \pi^*) \geq (1 - \alpha - \beta) \left(\left\lfloor \frac{H}{2} \right\rfloor - \frac{T_\tau}{2} \right) + \beta \left\lfloor \frac{H}{2} \right\rfloor.$$

1276 **Step 5: Bounding exploitability.** Finally, we will upper bound also the expected reward of the FH-MFG-NE policy π^* and
 1277 hence lower bound the exploitability. Our conclusion will be that π^* suffers from a non-vanishing exploitability for large H ,
 1278 as π^{br} becomes the best response policy after $H \gtrsim \log N$. In this step, we assume the probability space induced by all N
 1279 agents following FH-MFG-NE policy π^{br} .
 1280

1281 We have the definition

$$\begin{aligned}
 1282 \quad J_{P,R}^{H,N,(1)}(\pi^*, \pi^*, \dots, \pi^*) &= \mathbb{E} \left[\sum_{h=0}^{H-1} R(s_h^1, a_h^1, \hat{\mu}_h) \right] \\
 1283 &\leq (1 - \alpha - \beta) \mathbb{E} \left[\sum_{\text{odd } h=0}^{H-1} R_h^{1,g} \right] + (\alpha + \beta) \left\lfloor \frac{H}{2} \right\rfloor
 \end{aligned}$$

1284 This time, when h odd and $h > 2\tau$, it holds that $\mathbb{E}[R_h^g | h > 2\tau] = 1/2$ since π^* takes actions a_A, a_B with equal probability in
 1285 even steps, yielding $R_h^g = 1$ and $R_h^g = 0$ respectively almost surely. As before,
 1286

$$\begin{aligned}
 1287 \quad \mathbb{E} \left[\sum_{\substack{h \text{ odd} \\ 0 \leq h \leq H}} R_h^g \right] &= \mathbb{E} \left[\sum_{0 \leq h \leq \min\{2\tau, H\}} R_h^{1,g} + \sum_{\substack{h \text{ odd} \\ \min\{2\tau, H\} + 1 \leq h < H}} R_h^{1,g} \right] \\
 1288 &\leq \mathbb{E} \left[\min \left\{ \tau, \left\lfloor \frac{H}{2} \right\rfloor \right\} + \frac{1}{2} \left(\left\lfloor \frac{H}{2} \right\rfloor - \min \left\{ \tau, \left\lfloor \frac{H}{2} \right\rfloor \right\} \right) \right] \\
 1289 &= \frac{1}{2} \mathbb{E} \left[\left\lfloor \frac{H}{2} \right\rfloor + \min \left\{ \tau, \left\lfloor \frac{H}{2} \right\rfloor \right\} \right] \\
 1290 &\leq \frac{1}{2} \left\lfloor \frac{H}{2} \right\rfloor + \frac{1}{2} \mathbb{E}[\tau] = \frac{1}{2} \left\lfloor \frac{H}{2} \right\rfloor + \frac{1}{2} T_\tau.
 \end{aligned}$$

The statement of the theorem then follows by lower bounding the exploitability as follows:

$$\begin{aligned}
 & \mathcal{E}_{P,R}^{H,N,(1)}(\boldsymbol{\pi}^*, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) \\
 &= \max_{\boldsymbol{\pi}} J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) - J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}^*, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) \\
 &\geq J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}^{\text{br}}, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) - J_{P,R}^{H,N,(1)}(\boldsymbol{\pi}^*, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) \\
 &\geq (1 - \alpha - \beta) \left(\left\lfloor \frac{H}{2} \right\rfloor - \frac{T_\tau}{2} - \frac{1}{2} \left\lfloor \frac{H}{2} \right\rfloor - \frac{T_\tau}{2} \right) - \alpha \left\lfloor \frac{H}{2} \right\rfloor \\
 &\geq (1 - \alpha - \beta) \left(\frac{H}{4} - 24 - \frac{5}{9} \log_5 N \right) - \alpha \left\lfloor \frac{H}{2} \right\rfloor
 \end{aligned}$$

The above inequality implies that if $H \geq \log_2 N$, then

$$\begin{aligned}
 & \mathcal{E}_{P,R}^{H,N,(1)}(\boldsymbol{\pi}^*, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) \\
 &\geq (1 - \alpha - \beta) \left(\frac{1}{4} - \frac{5}{9 \log_2 5} \right) H - \alpha \frac{H}{2} - 24,
 \end{aligned}$$

which implies $\mathcal{E}_{P,R}^{H,N,(1)}(\boldsymbol{\pi}^*, \boldsymbol{\pi}^*, \dots, \boldsymbol{\pi}^*) \geq \Omega(H)$ by choosing α, β small constants as $\frac{1}{4} - \frac{5}{9 \log_2 5} > 0$.

A.4. Upper Bound for Stat-MFG: Extended Proof of Theorem 3.5

Let μ^*, π^* be a δ -Stat-MFG-NE. As before, the proof will proceed in three steps:

- **Step 1.** Bounding the expected deviation of the empirical population distribution from the mean-field distribution $\mathbb{E} [\|\widehat{\mu}_h - \mu^*\|_1]$ for any given policy $\boldsymbol{\pi}$.
- **Step 2.** Bounding difference of N agent value function $J_{P,R}^{\gamma,N,(i)}$ and the infinite player value function $V_{P,R}^\gamma$ in the stationary mean-field game setting.
- **Step 3.** Bounding the exploitability of an agent when each of N agents are playing the Stat-MFG-NE policy.

Step 1: Empirical distribution bound. We first analyze the deviation of the empirical population distribution $\widehat{\mu}_t$ over time from the stable distribution μ^* . For this, we state the following lemma and prove it using techniques similar to Corollary D.4 of (Yardim et al., 2023a).

Lemma A.11. *Assume that the conditions of Theorem 3.5 hold, and that $(\mu^*, \pi^*) \in \Delta_S$ is a Stat-MFG-NE. Furthermore, assume that the N agents follow policies $\{\pi^i\}_{i=1}^N$ in the N -Stat-MFG, define $\Delta_\pi := \frac{1}{N} \sum_i \|\pi^i - \pi^*\|_1$. Then, for any $t \geq 0$, we have*

$$\mathbb{E} [\|\mu^* - \widehat{\mu}_t\|_1] \leq \frac{tK_a \Delta_\pi}{2} + \frac{2(t+1)\sqrt{|S|}}{\sqrt{N}}.$$

Proof. \mathcal{F}_t as the σ -algebra generated by the states of agents $\{s_t^i\}$ at time t . For $\widehat{\mu}_0$, we have by definitions that

$$\begin{aligned}
 \mathbb{E} [\widehat{\mu}_0] &= \mathbb{E} \left[\frac{1}{N} \sum_i \mathbf{e}_{s_t^i} \right] = \mu^* \\
 \mathbb{E} [\|\widehat{\mu}_0 - \mu^*\|_2^2] &= \mathbb{E} \left[\frac{1}{N^2} \sum_i \left\| \left(\mathbf{e}_{s_t^i} - \mu^* \right) \right\|_2^2 \right] \leq \frac{4}{N}
 \end{aligned}$$

where the last line follows by independence. The two above imply $\mathbb{E} [\|\widehat{\mu}_0 - \mu^*\|_1] \leq \frac{2\sqrt{|S|}}{\sqrt{N}}$.

Next, we inductively calculate:

$$\begin{aligned} \mathbb{E}[\hat{\mu}_{t+1}|\mathcal{F}_t] &= \mathbb{E}\left[\frac{1}{N}\sum_{s'\in\mathcal{S}}\sum_{i=1}^N\mathbb{1}(s_{t+1}^i=s')\mathbf{e}_{s'}\middle|\mathcal{F}_t\right] \\ &= \sum_{s'\in\mathcal{S}}\mathbf{e}_{s'}\sum_{i=1}^N\frac{1}{N}\bar{P}(s'|s_t^i,\pi^i(s_t^i),\hat{\mu}_t), \end{aligned} \quad (10)$$

$$\begin{aligned} &\mathbb{E}[\|\hat{\mu}_{t+1}-\mathbb{E}[\hat{\mu}_{t+1}|\mathcal{F}_t]\|_2^2|\mathcal{F}_t] \\ &= \frac{1}{N^2}\sum_{i=1}^N\mathbb{E}[\|\mathbf{e}_{s_{t+1}^i}-\mathbb{E}[\mathbf{e}_{s_{t+1}^i}|\mathcal{F}_t]\|_2^2|\mathcal{F}_t] \leq \frac{4}{N}. \end{aligned} \quad (11)$$

We bound the ℓ_1 distance to the stable distribution as

$$\begin{aligned} &\mathbb{E}[\|\hat{\mu}_{t+1}-\mu^*\|_1|\mathcal{F}_t] \\ &\leq \underbrace{\mathbb{E}[\|\mathbb{E}[\hat{\mu}_{t+1}|\mathcal{F}_t]-\mu^*\|_1]}_{(\square)} + \underbrace{\mathbb{E}[\|\mathbb{E}[\hat{\mu}_{t+1}|\mathcal{F}_t]-\hat{\mu}_{t+1}\|_1|\mathcal{F}_t]}_{(\triangle)}. \end{aligned}$$

The two terms can be bounded separately using Inequalities (10) and (11).

$$\begin{aligned} (\triangle) &\leq \sqrt{|\mathcal{S}|}\mathbb{E}[\|\mathbb{E}[\hat{\mu}_{t+1}|\mathcal{F}_t]-\hat{\mu}_{t+1}\|_2|\mathcal{F}_t] \\ &\leq \sqrt{|\mathcal{S}|}\sqrt{\mathbb{E}[\|\mathbb{E}[\hat{\mu}_{t+1}|\mathcal{F}_t]-\hat{\mu}_{t+1}\|_2^2|\mathcal{F}_t]} \leq \frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}} \\ (\square) &= \left\|\sum_{s'\in\mathcal{S}}\mathbf{e}_{s'}\sum_{i=1}^N\frac{1}{N}\bar{P}(s'|s_t^i,\pi^i(s_t^i),\hat{\mu}_t)-\mu^*\right\|_1 \\ &= \left\|\sum_{s'\in\mathcal{S}}\mathbf{e}_{s'}\sum_{i=1}^N\frac{1}{N}\bar{P}(s'|s_t^i,\pi^i(s_t^i),\hat{\mu}_t)-\Gamma_{pop}(\pi^*,\mu^*)\right\|_1 \\ &\leq \left\|\sum_{i=1}^N\frac{1}{N}\bar{P}(\cdot|s_t^i,\pi^i(s_t^i),\hat{\mu}_t)-\sum_{i=1}^N\frac{1}{N}\bar{P}(\cdot|s_t^i,\pi^*(s_t^i),\hat{\mu}_t)\right\|_1 \\ &\quad + \left\|\sum_{s'\in\mathcal{S}}\hat{\mu}_t(s')\bar{P}(s'|s_t^i,\pi^i(s_t^i),\hat{\mu}_t)-\Gamma_{pop}(\pi^*,\mu^*)\right\|_1 \\ &\leq \frac{K_a}{2N}\sum_i\|\pi^*-\pi^i\|_1 + \|\Gamma_{pop}(\pi^*,\hat{\mu}_t)-\Gamma_{pop}(\pi^*,\mu^*)\|_1 \\ &\leq \frac{K_a\Delta\pi}{2} + \|\mu^*-\hat{\mu}_t\|_1 \end{aligned}$$

Hence, by the law of total expectation, we can conclude

$$\mathbb{E}[\|\mu^*-\hat{\mu}_{t+1}\|_1] \leq \mathbb{E}[\|\mu^*-\hat{\mu}_t\|_1] + \frac{K_a\Delta\pi}{2} + \frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}}$$

or inductively,

$$\mathbb{E}[\|\mu^*-\hat{\mu}_t\|_1] \leq \frac{tK_a\Delta\pi}{2} + \frac{2(t+1)\sqrt{|\mathcal{S}|}}{\sqrt{N}}.$$

□

Step 2: Bounding difference in value functions. Next, we bound the differences in the infinite-horizon

1430 **Lemma A.12.** *Suppose N -Stat-MFG agents follow the same sequence of policy π^* . Then for all i ,*

$$1431 \quad |J_{P,R}^{\gamma,N,(i)}(\pi^*, \dots, \pi^*) - V_{P,R}^\gamma(\mu^*, \pi^*)|$$

$$1432 \quad \leq \frac{\gamma}{1-\gamma} \left(L_\mu + \frac{L_s}{2} \right) \frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}}$$

1436 *Proof.* For ease of reading, in this proof expectations, probabilities, and laws of random variables will be denoted
 1437 $\mathbb{E}_\infty, \mathbb{P}_\infty, \mathcal{L}_\infty$ respectively over the infinite player finite horizon game and $\mathbb{E}_N, \mathbb{P}_N, \mathcal{L}_N$ respectively over the N -player game.
 1438 Due to symmetry in the N agent game, any permutation $\sigma : [N] \rightarrow [N]$ of agents does not change their distribution, that is
 1439 $\mathcal{L}_N(s_t^1, \dots, s_t^N) = \mathcal{L}_N(s_t^{\sigma(1)}, \dots, s_t^{\sigma(N)})$. We can then conclude that:

$$1441 \quad \mathbb{E}_N [R(s_t^1, a_t^1, \hat{\mu}_h)] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_N [R(s_t^i, a_t^i, \hat{\mu}_h)]$$

$$1442 \quad = \mathbb{E}_N \left[\sum_{s \in \mathcal{S}} \hat{\mu}_t(s) \bar{R}(s, \pi_t(s), \hat{\mu}_t) \right]$$

1446 Therefore, we by definition:

$$1448 \quad J_{P,R}^{\gamma,N,(1)}(\pi, \dots, \pi) = \mathbb{E}_N \left[\sum_{t=0}^{\infty} \sum_{s \in \mathcal{S}} \hat{\mu}_t(s) \bar{R}(s, \pi^*(s), \hat{\mu}_t) \right].$$

1451 Next, in the Stat-MFG, we have that for all $t \geq 0$,

$$1452 \quad \mathbb{P}_\infty(s_t = \cdot) = \mu^*,$$

$$1453 \quad \mathbb{P}_\infty(s_{t+1} = \cdot) = \sum_{s \in \mathcal{S}} \mathbb{P}_\infty(s_t = s) \mathbb{P}_\infty(s_t = \cdot | s_t = s)$$

$$1454 \quad = \Gamma_P(\mathbb{P}_\infty(s_t = s), \pi^*) = \mu^*,$$

1455 so by induction $\mathbb{P}_\infty(s_t = \cdot) = \mu^*$. Then we can conclude that

$$1459 \quad V_{P,R}^\gamma(\mu^*, \pi^*) = \mathbb{E}_\infty \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi^*(s_t), \mu_t) \right]$$

$$1460 \quad = \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \mu^*(s) R(s, \pi^*(s), \mu^*),$$

1465 by a simple application of the dominated convergence theorem. We next bound the differences in truncated expect reward
 1466 until some time $T > 0$:

$$1467 \quad \left| \mathbb{E}_N \left[\sum_{t=0}^T \gamma^t \sum_{s \in \mathcal{S}} \hat{\mu}_t(s) \bar{R}(s, \pi^*(s), \hat{\mu}_t) \right] \right.$$

$$1468 \quad \left. - \sum_{t=0}^T \gamma^t \sum_{s \in \mathcal{S}} \mu_t(s) R(s, \pi^*(s), \mu_t) \right|$$

$$1470 \quad \leq \mathbb{E}_N \left[\sum_{t=0}^T \gamma^t \left| \sum_{s \in \mathcal{S}} (\hat{\mu}_t(s) \bar{R}(s, \pi^*(s), \hat{\mu}_t) - \mu^*(s) R(s, \pi^*(s), \mu^*)) \right| \right]$$

$$1471 \quad \leq \mathbb{E}_N \left[\sum_{t=0}^T \gamma^t \left(\frac{L_s}{2} \|\mu^* - \hat{\mu}_t\|_1 + L_\mu \|\mu^* - \hat{\mu}_t\|_1 \right) \right]$$

$$1472 \quad \leq \sum_{t=0}^T \gamma^t \left(L_\mu + \frac{L_s}{2} \right) \mathbb{E}_N [\|\mu^* - \hat{\mu}_t\|_1]$$

$$1473 \quad \leq \frac{1}{(1-\gamma)^2} \left(L_\mu + \frac{L_s}{2} \right) \frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}}$$

1485 Taking $T \rightarrow \infty$ and applying once again the dominated convergence theorem the result is obtained. \square

1486

1487

1488 **Step 3: Bounding difference in policy deviation.** Finally, to conclude the proof of the main theorem of this section, we
 1489 will prove that the improvement in expectation due to single-sided policy changes are at most of order $\mathcal{O}\left(\delta + \frac{1}{\sqrt{N}}\right)$.

1490

1491 **Lemma A.13.** *Suppose we have two policy sequences $\pi^*, \pi \in \Pi$ and $\mu^* \in \Delta_{\mathcal{S}}$ such that $\Gamma_P(\mu^*, \pi^*) = \mu^*$ and $\Gamma_P(\cdot, \pi^*)$
 1492 is non-expansive. Then,*

1493

$$\begin{aligned} & \left| J_{P,R}^{\gamma, N, (1)}(\pi', \pi^*, \dots, \pi^*) - V_{P,R}^{\gamma}(\mu^*, \pi') \right| \\ & \leq \sum_{t=0}^{\infty} \gamma^t \left(L_{\mu} \mathbb{E} [\|\hat{\mu}_t - \mu_t^{\pi}\|_1] + K_{\mu} \sum_{t'=0}^{t-1} \mathbb{E} [\|\hat{\mu}_{t'} - \mu_{t'}^{\pi}\|_1] \right) \\ & \leq \left(\frac{K_a}{2N} + \frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}} \right) \frac{L_{\mu}/2 + K_{\mu}}{(1-\gamma)^3} \end{aligned}$$

1499

1500

1501 *Proof.* For the truncated game T , it still holds by the derivation in the FH-MFG that:

1502

$$\begin{aligned} & \left| \mathbb{E}_N [R(s_t^1, a_t^1, \hat{\mu}_t)] - \mathbb{E}_{\infty} [R(s_t, a_t, \mu_t^{\pi})] \right| \\ & \leq \frac{L_{\mu}}{2} \mathbb{E}_N [\|\mu_t^{\pi} - \hat{\mu}_t\|_1] + K_{\mu} \sum_{t'=0}^{t-1} \mathbb{E}_N [\|\mu_{t'}^{\pi} - \hat{\mu}_{t'}\|_1]. \end{aligned}$$

1507

1508 We take the limit $T \rightarrow \infty$ and apply the dominated convergence theorem to obtain the state bound, also noting that
 1509 $1/2 \cdot \sum_t (t+1)(t+2)\gamma^t \leq \frac{1}{(1-\gamma)^3}$. \square

1510

1511

1512 **Conclusion and Statement of the Result.** Finally, if μ^*, π^* is a δ -Stat-MFG-NE, by definition we have that: By definition
 1513 of the Stat-MFG-NE, we have:

1514

$$\delta \geq \mathcal{E}_{P,R}^H(\pi_{\delta}) = \max_{\pi' \in \Pi} V_{P,R}^{\gamma}(\mu^*, \pi') - V_{P,R}^{\gamma}(\mu^*, \pi^*)$$

1517

1518 Then using the two bounds from Steps 2,3 and the fact that π^* δ -optimal with respect to μ^* :

1519

$$\begin{aligned} & \max_{\pi' \in \Pi} J_{P,R}^{H, N, (1)}(\pi', \pi^*, \dots, \pi^*) - J_{P,R}^{H, N, (1)}(\pi^*, \pi^*, \dots, \pi^*) \\ & \leq 2\delta + \left(\frac{K_a}{2N} + \frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}} \right) \frac{L_{\mu}/2 + K_{\mu}}{(1-\gamma)^3} + \frac{L_{\mu} + L_s/2}{(1-\gamma)^2} \left(\frac{2\sqrt{|\mathcal{S}|}}{\sqrt{N}} \right) \end{aligned}$$

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

A.5. Lower Bound for Stat-MFG: Extended Proof of Theorem 3.6

Similar to the finite horizon case, we define constructively the counter-example: the idea and the nature of the counter-example remain the same. However, minor details of the construction are modified, as it will not hold immediately that all agents are on states $\{s_{\text{Left}}, s_{\text{Right}}\}$ on even times t , and that the Stat-MFG-NE is unique as before.

Defining the Stat-MFG. We use the same definitions for $\mathcal{S}, \mathcal{A}, \mathbf{g}, \mathbf{h}, \omega_{\varepsilon}$ as in the FH-MFG case. Define the convenience functions Q_L, Q_R as

$$\begin{aligned} Q_L(\mu) & := \frac{\mu(s_{\text{LA}}) + \mu(s_{\text{LB}})}{\max\{\mu(s_{\text{LA}}) + \mu(s_{\text{LB}}) + \mu(s_{\text{RA}}) + \mu(s_{\text{RB}}), 4/9\}}, \\ Q_R(\mu) & := \frac{\mu(s_{\text{RA}}) + \mu(s_{\text{RB}})}{\max\{\mu(s_{\text{LA}}) + \mu(s_{\text{LB}}) + \mu(s_{\text{RA}}) + \mu(s_{\text{RB}}), 4/9\}}. \end{aligned}$$

1540 We define the transition probabilities:

1541
1542
1543
1544
1545
1546
1547
1548

$$\text{If } s \in \{s_{LA}, s_{LB}, s_{RA}, s_{RB}\}, \forall \mu, a :$$

$$P(s'|s, a, \mu) = \begin{cases} \omega_\epsilon(Q_L(\mu)), & \text{if } s' = s_{Right}, s \in \{s_{LA}, s_{LB}\} \\ \omega_\epsilon(Q_R(\mu)), & \text{if } s' = s_{Left}, s \in \{s_{LA}, s_{LB}\} \\ \omega_\epsilon(Q_L(\mu)), & \text{if } s' = s_{Right}, s \in \{s_{RA}, s_{RB}\} \\ \omega_\epsilon(Q_R(\mu)), & \text{if } s' = s_{Left}, s \in \{s_{RA}, s_{RB}\} \end{cases},$$

1549 and define $P(s_{Left}, a, \mu), P(s_{Right}, a, \mu)$ as before. With previous Lipschitz continuity results, it follows that $P \in \mathcal{P}_{9/s\epsilon}$.

1551 Similarly, we modify the reward function R as follows:

1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567

$$\begin{aligned} R(s_{Left}, a_A, \mu) &= R(s_{Left}, a_B, \mu) = 0, \\ R(s_{Right}, a_A, \mu) &= R(s_{Right}, a_B, \mu) = 0, \\ \begin{pmatrix} R(s_{LA}, a_A, \mu) \\ R(s_{LB}, a_A, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(Q_L(\mu), Q_R(\mu)) + \alpha \mathbf{h}(\mu(s_{LA}), \mu(s_{LB})) \\ \begin{pmatrix} R(s_{LA}, a_B, \mu) \\ R(s_{LB}, a_B, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(Q_L(\mu), Q_R(\mu)) + \mathbf{h}(\mu(s_{LA}), \mu(s_{LB})) \\ &\quad + \beta \mathbf{1} \\ \begin{pmatrix} R(s_{RA}, a_A, \mu) \\ R(s_{RB}, a_A, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(Q_R(\mu), Q_L(\mu)) + \alpha \mathbf{h}(\mu(s_{RA}), \mu(s_{RB})) \\ \begin{pmatrix} R(s_{RA}, a_B, \mu) \\ R(s_{RB}, a_B, \mu) \end{pmatrix} &= (1 - \alpha - \beta) \mathbf{g}(Q_R(\mu), Q_L(\mu)) + \alpha \mathbf{h}(\mu(s_{RA}), \mu(s_{RB})) \\ &\quad + \beta \mathbf{1}, \end{aligned}$$

1568 simple computation shows that $R \in \mathcal{R}_3$. In this proof, unlike the N -FH-SAG case, α will be chosen as a function of N ,
1569 namely $\alpha = \mathcal{O}(e^{-N})$.

1571 **Step 1: Solution of the Stat-MFG.** We solve the infinite agent game: let μ^*, π^* be an Stat-MFG-NE. By simple computation,
1572 one can see that for any stationary distribution μ^* of the game, probability must be distributed equally between groups of
1573 states $\{s_{Left}, s_{Right}\}$ and $\{s_{LA}, s_{LB}, s_{RA}, s_{RB}\}$, that is,

1574
1575
1576
1577

$$\begin{aligned} \mu^*(s_{Left}) + \mu^*(s_{Right}) &= 1/2, \\ \mu^*(s_{LA}) + \mu^*(s_{LB}) + \mu^*(s_{RA}) + \mu^*(s_{RB}) &= 1/2. \end{aligned}$$

1578 It holds by the stationarity equation $\Gamma_P(\mu^*, \pi^*) = \pi^*$ that

1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590

$$\begin{aligned} \mu^*(s_{Left}) &= \mu^*(s_{LA}) + \mu^*(s_{LB}), \\ \mu^*(s_{Right}) &= \mu^*(s_{RA}) + \mu^*(s_{RB}), \\ \mu^*(s_{Left}) &= \sum_{s \in \mathcal{S}} \mu^*(s) \pi^*(a|s) P(s_{Left}|s, a, \mu^*) \\ &= P(s_{Left}|s_{LA}, a_A, \mu^*), \\ \mu^*(s_{Right}) &= \sum_{s \in \mathcal{S}} \mu^*(s) \pi^*(a|s) P(s_{Right}|s, a, \mu^*) \\ &= P(s_{Right}|s_{LA}, a_A, \mu^*), \end{aligned}$$

1591 as $P(s_{Right}|s, a, \mu^*) = P(s_{Right}|s, a, \mu^*)$ and similarly $P(s_{Left}|s, a, \mu^*) = P(s_{Left}|s, a, \mu^*)$ for any $s \in$
1592 $\{s_{LA}, s_{LB}, s_{RA}, s_{RB}\}, a \in \mathcal{A}$. If $\mu^*(s_{Left}) > 1/4$, then by definition $P(s_{Left}|s_{LA}, a_A, \mu^*) < 1/4$, and similarly if
1593 $\mu^*(s_{Left}) < 1/4$, then by definition $P(s_{Left}|s_{LA}, a_A, \mu^*) > 1/4$. So it must be the case that $\mu^*(s_{Left}) = \mu^*(s_{Right}) = 1/4$. Then
1594

1595 the unique Stat-MFG-NE must be

1596

1597

1598

1599

1600

1601

$$\pi^*(a|s) := \begin{cases} 1, & \text{if } a = a_B, s \in \{s_{LA}, s_{LB}, s_{RA}, s_{RB}\} \\ \frac{1}{2}, & \text{if } s \in \{s_{Left}, s_{Right}\} \\ 0, & \text{if } a = a_A, s \in \{s_{LA}, s_{LB}, s_{RA}, s_{RB}\}, \end{cases}$$

$$\mu^*(s_{RA}) = \mu^*(s_{LA}) = \mu^*(s_{RB}) = \mu^*(s_{LB}) = 1/8,$$

1602

1603

1604

as otherwise the action $\arg \min_{a \in \mathcal{A}} \pi^*(a|s_{Right})$ will be a better response in state s_{Right} and the action $\arg \min_{a \in \mathcal{A}} \pi^*(a|s_{Left})$ will be optimal in state s_{Right} .

1605

1606

1607

1608

Step 2: Expected population deviation in N -Stat-SAG. We fix $1/2\varepsilon = 3$, define the random variable $\bar{N} := N(\hat{\mu}_0(s_{Right}) + \hat{\mu}_0(s_{Left}))$. We will analyze the population under the event $\bar{N} := \{|\bar{N}/N - 1/2| \leq 1/18\}$, which holds with probability $\Omega(1 - e^{-N^2})$ by the Hoeffding inequality. Under the event \bar{E} , it holds that $\hat{\mu}_t(s_{LA}) + \hat{\mu}_t(s_{LA}) + \hat{\mu}_t(s_{LA}) + \hat{\mu}_t(s_{LA}) > 4/9$ almost surely at all t .

1609

1610

1611

Fix $N_0 \in \mathbb{N}_{>0}$ such that $|N_0/N - 1/2| \leq 1/18$, in this step we will condition on $E_0 := \{\bar{N} := N_0\}$. Once again define the random process X_m for $m \in \mathbb{N}_{\geq 0}$ such that

1612

1613

1614

1615

$$X_m := \begin{cases} \frac{\hat{\mu}_{2m}(s_{Left})}{\hat{\mu}_{2m}(s_{Left}) + \hat{\mu}_{2m}(s_{Right})}, & \text{if } m \text{ odd} \\ \frac{\hat{\mu}_{2m}(s_{Right})}{\hat{\mu}_{2m}(s_{Left}) + \hat{\mu}_{2m}(s_{Right})}, & \text{if } m \text{ even} \end{cases}$$

1616

1617

1618

with the modification at odd m necessary because of the difference in dynamics P (oscillating between s_{Left}, s_{Right}) from the FH-SAG case. It still holds that X_m is Markovian, and given X_m we have $N_0 X_{m+1} \sim \text{Binom}(N_0, \omega_\varepsilon(X_m))$. As before, X_m is independent from the policies of agents.

1619

1620

1621

1622

1623

1624

Define $K := \lfloor \log_2 \sqrt{N_0} \rfloor$, $\mathcal{G} := \{k/N_0 : k = 0, \dots, N_0\}$, $\mathcal{G}_* := \{0, 1\} \subset \mathcal{G}$ and the level sets once again as

$$\mathcal{G}_{-1} := \mathcal{G}, \quad \mathcal{G}_k := \left\{ x \in \mathcal{G} : \left| x - \frac{1}{2} \right| \geq \frac{2^k}{2\sqrt{N_0}} \right\} \text{ when } k \leq K,$$

$$\mathcal{G}_{K+1} := \mathcal{G}_*.$$

1625

1626

1627

As before, using the Markov property, Hoeffding, and the fact that $|\omega_\varepsilon(x) - 1/2| \geq 1/2\varepsilon|x - 1/2|$ we obtain $\forall k \in 0, \dots, K-1$, $\forall m$ that

1628

1629

1630

1631

$$\mathbb{P}[X_{m+1} \in \mathcal{G}_0 | X_m \in \mathcal{G}_{-1}, E_0] \geq 1/20$$

$$\mathbb{P}[X_{m+1} \in \mathcal{G}_{k+1} | X_m \in \mathcal{G}_k, E_0] \geq \alpha_k := 1 - 2 \exp \left\{ -\frac{1}{8} 4^{k+1} \right\},$$

1632

1633

1634

1635

hence from the analysis before we have the lower bound

$$\mathbb{E}[|X_m - 1/2| | E_0] \geq C_1 \min \left\{ \frac{2^m}{\sqrt{N_0}}, 1 \right\},$$

1636

1637

for some absolute constant $C_2 > 0$.

1638

1639

Step 3. Exploitability lower bound. As in the case of FH-MFG, the ergodic optimal policy is given by

1640

1641

1642

1643

1644

$$\bar{\pi}(a|s) = \begin{cases} 1, & \text{if } s = s_{Left}, a = a_A \\ 1, & \text{if } s = s_{Right}, a = a_A \\ 1, & \text{if } s \notin \{s_{Left}, s_{Right}\}, a = a_B \\ 0, & \text{otherwise} \end{cases}$$

1645

1646

1647

1648

1649

We define the shorthand functions

$$\mathcal{S}^* := \{s_{Left}, s_{Right}\}, \quad Q(\mu) := (Q_L(\mu), Q_R(\mu)),$$

$$Q_{\min}(\mu) := \min\{Q_L(\mu), Q_R(\mu)\}, \quad Q_{\max} := \max\{Q_L(\mu), Q_R(\mu)\}.$$

1650 We condition on $E_{\mathcal{S}^*} := \{s_0^1 \in \mathcal{S}^*\}$, that is the first agent starts from states $\{s_{\text{Left}}, s_{\text{Right}}\}$, the analysis will be similar under
 1651 event $E_{\mathcal{S}^*}^c$. As in the case of FH-MFG, due to permutation invariance, it holds for any odd t and $\mu \in \{\mu' \in \Delta_{\mathcal{S}^*} : N_0 \mu' \in$
 1652 $\mathbb{N}_{>0}^2\}$ that

$$1653 \quad \mathbb{P}[s_t^1 \in \{s_{\text{LA}}, s_{\text{LB}}\} | E_0, E_{\mathcal{S}^*}, Q(\hat{\mu}_t) = \mu] = Q_L(\mu)$$

$$1654 \quad \mathbb{P}[s_t^1 \in \{s_{\text{RA}}, s_{\text{RB}}\} | E_0, E_{\mathcal{S}^*}, Q(\hat{\mu}_t) = \mu] = Q_R(\mu),$$

1655 therefore expressing the error component due to \mathbf{g} as $R_t^{1,\mathbf{g}}$ and expressing some repeating conditionals as \bullet :

$$1656 \quad \bar{G}_t^\mu := \mathbb{E} \left[R_t^{1,\mathbf{g}} \middle| E_0, E_{\mathcal{S}^*}, Q(\hat{\mu}_t) = \mu, a_t^1 \sim \bar{\pi}(s_t^1), \begin{smallmatrix} a_i^i \sim \pi^*(s_t^i), \\ \text{when } i \neq 1 \end{smallmatrix} \right]$$

$$1657 \quad = \sum_{s \in \mathcal{S}^*} \mathbb{P}[s_t^1 = s | Q(\hat{\mu}_t) = \mu, \bullet] \mathbb{E}[R_t^{1,\mathbf{g}} | s_t^1 = s, Q(\hat{\mu}_t) = \mu, \bullet]$$

$$1658 \quad = \frac{Q_{\max}(\mu)}{Q_{\max}(\mu)} Q_{\max}(\mu) + \frac{Q_{\min}(\mu)}{Q_{\max}(\mu)} Q_{\min}(\mu).$$

1659 Similarly, since $\pi^*(a|s) = 1/2$ for any $s \in \mathcal{S}^*$, it holds that

$$1660 \quad G_t^\mu := \mathbb{E} \left[R_t^{1,\mathbf{g}} \middle| E_0, E_{\mathcal{S}^*}, Q(\hat{\mu}_t) = \mu, \begin{smallmatrix} a_i^i \sim \pi^*(s_t^i), \\ \forall i \end{smallmatrix} \right]$$

$$1661 \quad = \frac{1}{2} \frac{Q_{\min}(\mu)}{Q_{\max}(\mu)} + \frac{1}{2} \frac{Q_{\max}(\mu)}{Q_{\max}(\mu)}.$$

1662 Therefore, given the population distribution between $s_{\text{LA}}, s_{\text{LB}}$ and $s_{\text{RA}}, s_{\text{RB}}$, the expected difference in rewards for the two
 1663 policies is

$$1664 \quad \bar{G}_t^\mu - G_t^\mu = \left(Q_{\max}(\mu) - \frac{1}{2} \right) + \left(Q_{\min}(\mu) - \frac{1}{2} \right) \frac{Q_{\min}(\mu)}{Q_{\max}(\mu)}$$

$$1665 \quad = \left(Q_{\max}(\mu) - \frac{1}{2} \right) + \left(\frac{1}{2} - Q_{\max}(\mu) \right) \frac{Q_{\min}(\mu)}{Q_{\max}(\mu)}$$

$$1666 \quad = \left(Q_{\max}(\mu) - \frac{1}{2} \right) \left(1 - \frac{Q_{\min}(\mu)}{Q_{\max}(\mu)} \right)$$

$$1667 \quad \geq 2 \left(Q_{\max}(\mu) - \frac{1}{2} \right)^2.$$

1668 Therefore from above, we conclude that

$$1669 \quad \mathbb{E}[\bar{G}_t^{\hat{\mu}_t} - G_t^{\hat{\mu}_t} | E_0] \geq \mathbb{E}[2|X_{t-\frac{1}{2}} - 1/2|^2 | E_0, E_{\mathcal{S}^*}] \geq 2C_1^2 \min \left\{ \frac{2^t}{2N_0}, 1 \right\}.$$

1670 Using the lower bound above, the conditional expected difference in discounted total reward is

$$1671 \quad \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^1, a_t^1, \hat{\mu}_t) \middle| E_0, E_{\mathcal{S}^*}, a_t^1 \sim \bar{\pi}(s_t^1), \begin{smallmatrix} a_i^i \sim \pi^*(s_t^i), \\ \text{when } i \neq 1 \end{smallmatrix} \right]$$

$$1672 \quad - \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^1, a_t^1, \hat{\mu}_t) \middle| E_0, E_{\mathcal{S}^*}, \begin{smallmatrix} a_i^i \sim \pi^*(s_t^i), \\ \forall i \end{smallmatrix} \right]$$

$$1673 \quad \geq (1 - \alpha - \beta) \sum_{k=0}^{\infty} 2C_1^2 \gamma^{2k+1} \min \left\{ \frac{2^{2k}}{N_0}, 1 \right\} - \frac{2\alpha}{1 - \gamma}$$

$$1674 \quad \geq \frac{C_2}{N_0} \sum_{k=0}^{\lceil \log_4 N_0 \rceil} (4\gamma^2)^k + \frac{C_3}{N_0} \sum_{k=\lceil \log_4 N_0 \rceil}^{\infty} \gamma^{2k} - \frac{2\alpha}{1 - \gamma}$$

$$1675 \quad \geq \frac{C_4((4\gamma^2)^{\log_4 N_0} - 1)}{N_0} + C_5 \frac{(\gamma^2)^{\log_4 N_0} N_0^{-1}}{1 - \gamma^2} - \frac{2\alpha}{1 - \gamma}$$

$$1676 \quad \geq C_6 N_0^{\log_2 \gamma} + C_7 \frac{N_0^{\log_2 \gamma - 1}}{1 - \gamma} - \frac{2\alpha}{1 - \gamma}.$$

1705 Taking expectation over N_0 (using $\mathbb{E}[\bar{N}|E^*] = N/2$ and Jensen's):

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^1, a_t^1, \hat{\mu}_t) | E^*, E_{S^*}, a_t^1 \sim \bar{\pi}(s_t^1), \begin{smallmatrix} a_t^i \sim \pi^*(s_t^i), \\ \text{when } i \neq 1 \end{smallmatrix} \right] \\
 & - \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^1, a_t^1, \hat{\mu}_t) | E^*, E_{S^*}, a_t^i \sim \pi^*(s_t^i), \forall i \right] \\
 & \geq C_6 N_0^{\log_2 \gamma} + C_7 \frac{N_0^{\log_2 \gamma - 2}}{1 - \gamma} - \frac{2\alpha}{1 - \gamma}
 \end{aligned}$$

1710 While the analysis above assumes event E_{S^*} , the same analysis lower bound follows with a shift between even and odd
 1711 steps when $s_0^1 \notin S^*$, hence

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^1, a_t^1, \hat{\mu}_t) | E^*, a_t^1 \sim \bar{\pi}(s_t^1), \begin{smallmatrix} a_t^i \sim \pi^*(s_t^i), \\ \text{when } i \neq 1 \end{smallmatrix} \right] \\
 & - \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^1, a_t^1, \hat{\mu}_t) | E^*, a_t^i \sim \pi^*(s_t^i), \forall i \right] \\
 & \geq C_6 N_0^{\log_2 \gamma} + C_7 \frac{N_0^{\log_2 \gamma - 2}}{1 - \gamma} - \frac{2\alpha}{1 - \gamma}
 \end{aligned}$$

1727 Finally, we conclude the proof with the observation

$$\begin{aligned}
 & \max_{\pi} J_{P,R}^{\gamma, N, (1)}(\pi, \pi^*, \dots, \pi^*) - J_{P,R}^{H, N, (1)}(\pi^*, \pi^*, \dots, \pi^*) \\
 & \geq J_{P,R}^{\gamma, N, (1)}(\bar{\pi}, \pi^*, \dots, \pi^*) - J_{P,R}^{H, N, (1)}(\pi^*, \pi^*, \dots, \pi^*) \\
 & \geq C_6 N_0^{\log_2 \gamma} + C_7 \frac{N_0^{\log_2 \gamma - 2}}{1 - \gamma} - \frac{2\alpha}{1 - \gamma} - (1 - \gamma)^{-1} \mathbb{P}[\bar{E}^c],
 \end{aligned}$$

1735 where $\mathbb{P}[\bar{E}^c] = O(e^{-N^2})$ and we pick $\alpha = O(e^{-N})$.

1738 B. Intractability Results

1739 B.1. Fundamentals of PPAD

1741 We first introduce standard definitions and tools, mostly taken from (Daskalakis et al., 2009; Goldberg, 2011; Papadimitriou,
 1742 1994).

1744 **Notations.** For a finite set Σ , we denote by Σ^n the set of tuples n elements from Σ , and by $\Sigma^* = \bigcup_{n \geq 0} \Sigma^n$ the set of finite
 1745 sequences of elements of Σ . For any $\alpha \in \Sigma$, let $\alpha^n \in \Sigma^n$ denote the n -tuple $(\underbrace{\alpha, \dots, \alpha}_{n \text{ times}})$. For $x \in \Sigma^*$, by $|x|$ we denote the
 1746 length of the sequence x . Finally, the following function will be useful, defined for any $\alpha > 0$:

$$\begin{aligned}
 & u_{\alpha} : \mathbb{R} \rightarrow [0, \alpha] \\
 & u_{\alpha}(x) := \max\{0, \min\{\alpha, x\}\} = \begin{cases} \alpha, & \text{if } x \geq \alpha, \\ x, & \text{if } 0 \leq x \leq \alpha, \\ 0, & \text{if } x \leq 0. \end{cases}
 \end{aligned}$$

1755 We define a search problem \mathcal{S} on alphabet Σ as a relation from a set $\mathcal{I}_{\mathcal{S}} \subset \Sigma^*$ to Σ^* such that for all $x \in \mathcal{I}_{\mathcal{S}}$, the image of x
 1756 under \mathcal{S} satisfies $\mathcal{S}_x \subset \Sigma^{|x|^k}$ for some $k \in \mathbb{N}_{>0}$, and given $y \in \Sigma^{|x|^k}$ whether $y \in \mathcal{S}_x$ is decidable in polynomial time.

1757 Intuitively speaking, PPAD is the complexity class of search problems that can be shown to always have a solution using a
 1758 ‘‘parity argument’’ on a directed graph. The simplest complete example (the example that defines the problem class) of PPAD
 1759

problems is the computational problem END-OF-THE-LINE. The problem, formally defined below, can be summarized as such: given a directed graph where each node has in-degree and out-degree at most one and given a node that is a source in this graph (i.e., no incoming edge but one outgoing edge), find another node that is a sink or a source. Such a node can be always shown to exist using a simple parity argument.

Definition B.1 (END-OF-THE-LINE (Daskalakis et al., 2009)). The computational problem END-OF-THE-LINE is defined as follows: given two binary circuits S, P each with n input bits and n output bits such that $P(0^n) = 0^n \neq S(s^n)$, find an input $x \in \{0, 1\}^n$ such that $P(S(x)) \neq x$ or $S(P(x)) \neq x \neq 0^n$.

The obvious solution to the above is to follow the graph node by node using the given circuits until we reach a sink: however, this can take exponential time as the graph size can be exponential in the bit descriptions of the circuits. It is believed that END-OF-THE-LINE is difficult (Goldberg, 2011), that there is no efficient way to use the bit descriptions of the circuits S, P to find another node with degree 1.

B.2. Proof of Intractability of Stat-MFG

We reduce any ε -GCIRCUIT problem to the problem ε -STATDIST for some simple transition function $P \in \mathcal{P}^{\text{Sim}}$.

Let $(\mathcal{V}, \mathcal{G})$ be a generalized circuit to be reduced to a stable distribution computation problem. Let $V = |\mathcal{V}| \geq 1$. We will define a game that has at most $V + 1$ states and $|\mathcal{A}| = 1$ actions, that is, agent policy will not have significance, and it will suffice to determine simple transition probabilities $P(s'|s, \mu)$ for all $s, s' \in \mathcal{S}, \mu \in \Delta_{\mathcal{S}}$.

The proposed system will have a base state $s_{\text{base}} \in \mathcal{S}$ and 1 additional state s_v associated with the gate whose output is $v \in \mathcal{V}$. Our construction will be sparse: only transition probabilities in between states associated with a gate and s_{base} will take positive values. We define the useful constants $\theta := \frac{1}{8V}, B := \frac{1}{4}$.

Given an (approximately) stable distribution μ^* of P , for each vertex v we will read the satisfying assignment for the ε -GCIRCUIT problem by the value $u_1(\theta^{-1}\mu^*(s_v))$. For each possible gate, we define the following gadgets.

Binary assignment gadget. For a gate of the form $G_{\leftarrow}(\zeta||v)$, we will add one state s_v such that

$$\text{If } \zeta = 1 : \begin{cases} P(s_{\text{base}}|s_v, \mu) = 1, \\ P(s_v|s_v, \mu) = 0, \\ P(s_v|s_{\text{base}}, \mu) = \frac{\theta}{\max\{B, \mu(s_{\text{base}})\}} \end{cases}$$

$$\text{If } \zeta = 0 : \begin{cases} P(s_{\text{base}}|s_v, \mu) = 1, \\ P(s_v|s_v, \mu) = 0, \\ P(s_v|s_{\text{base}}, \mu) = 0 \end{cases}$$

Weighted addition gadget. Next, we implement the addition gadget $G_{\times,+}(\alpha, \beta|v_1, v_2|v)$ for $\alpha, \beta \in [-1, 1]$. In this case, we also add one state s_v to the game, and define the transition probabilities:

$$P(s_{\text{base}}|s_v, \mu) = 1,$$

$$P(s_v|s_v, \mu) = 0,$$

$$P(s_v|s_{\text{base}}, \mu) = \frac{u_{\theta}(\alpha u_{\theta}(\mu(v_1)) + \beta u_{\theta}(\mu(v_2)))}{\max\{B, \mu(s_{\text{base}})\}}$$

Brittle comparison gadget. For the comparison gate $G_{<}(|v_1, v_1|v)$, we also add one state s_v to the game. Define the function $p_{\delta} : [-1, 1] \rightarrow [0, 1]$

$$p_{\delta}(x, y) := u_1 \left(\frac{1}{2} + \delta^{-1}(x - y) \right),$$

1815 for any $\delta > 0$. In particular, if $x \geq y + \delta$, then $p_\delta(x, y) = 1$, and if $x \leq y - \delta$, then $p_\delta(x, y) = 0$. We define the probability
1816 transitions to and from s_v as

$$1817$$

$$1818 \quad P(s_v | s_{\text{base}}, \mu) = \frac{\theta p_{8\varepsilon}(\theta^{-1}u_\theta(\mu(s_1)), \theta^{-1}u_\theta(\mu(s_2)))}{\max\{B, \mu(s_{\text{base}})\}},$$

$$1819$$

$$1820 \quad P(s_v | s_v, \mu) = 0,$$

$$1821 \quad P(s_{\text{base}} | s_v, \mu) = 1.$$

$$1822$$

1823 Finally, after all s_v have been added, we complete the definition of P by setting

$$1824 \quad P(s_{\text{base}} | s_{\text{base}}, \mu) = 1 - \sum_{s' \in \mathcal{S}} P(s' | s_{\text{base}}, \mu).$$

1825 We first verify that the above assignment is a valid transition probability matrix for any $\mu \in \Delta_{\mathcal{S}}$. It is clear from definitions
1826 that for any $\mu, s \neq s_{\text{base}}$, $P(\cdot | s, \mu)$ is a valid probability distribution as long as $8\varepsilon < 1$. Moreover, for any $s \neq s_{\text{base}}$, it holds
1827 that $0 \leq P(s | s_{\text{base}}, \mu) \leq \frac{\theta}{B} < 1$, and it also holds that

$$1828 \quad P(s_{\text{base}} | s_{\text{base}}, \mu) = 1 - \sum_{s' \in \mathcal{S}} P(s' | s_{\text{base}}, \mu) \geq 1 - \frac{V\theta}{B} \geq 0$$

1829 so $P(\cdot | s_{\text{base}}, \mu)$ is a valid probability transition matrix. Finally, the defined transition probability function P is Lipschitz in
1830 the components of μ , and P can be defined as a composition of simple functions, hence $P \in \mathcal{P}^{\text{Sim}}$. Finally, in this defined
1831 MFG, it holds that $V + 1 = |\mathcal{S}|$, since for each gate in the generalized circuit we defined one additional state.

1832 **Error propagation.** We finally analyze the error propagation of the stationary distribution problem in terms of the
1833 generalized circuit. Without loss of generality we assume $\varepsilon < \frac{1}{8}$. First, for any solution of the ε -STATDIST problem μ^* ,
1834 whenever $\varepsilon < \frac{1}{8}$, it must hold that:

$$1835 \quad \left| \mu^*(s_{\text{base}}) - \sum_{s' \in \mathcal{S}} \mu^*(s) P(s_{\text{base}} | s, \mu^*) \right| \leq \frac{1}{8|\mathcal{S}|},$$

1836 hence (using $V < |\mathcal{S}|$) we have the lower bound on $\mu^*(s_{\text{base}})$ given by:

$$1837 \quad \mu^*(s_{\text{base}}) \geq \sum_{s \in \mathcal{S}} \mu^*(s) P(s_{\text{base}} | s, \mu^*) - \frac{1}{8V}$$

$$1838 \quad \geq \mu^*(s_{\text{base}}) P(s_{\text{base}} | s_{\text{base}}, \mu^*) + \sum_{s \neq s_{\text{base}}} \mu^*(s) P(s_{\text{base}} | s, \mu^*) - \frac{1}{8V}$$

$$1839 \quad \geq \mu^*(s_{\text{base}}) \left(1 - \frac{V\theta}{B}\right) + \sum_{s \neq s_{\text{base}}} \mu^*(s) - \frac{1}{8V}$$

$$1840 \quad \geq \mu^*(s_{\text{base}}) \left(1 - \frac{V\theta}{B}\right) + (1 - \mu^*(s_{\text{base}})) - \frac{1}{8V}$$

$$1841 \quad \implies \mu^*(s_{\text{base}}) \geq \frac{1 - \frac{1}{8V}}{1 + \frac{V\theta}{B}} \geq B = \frac{1}{4}.$$

1842 We will show that a solution of the ε -STATDIST can be converted into a ε' -satisfying assignment

$$1843 \quad v \rightarrow u_1 \left(\frac{\mu^*(s_v)}{\theta} \right),$$

1844 for some appropriate ε' to be defined later.

1870 **Case 1: Binary assignment error.** First, assume $G_{\leftarrow}(\zeta|v) \in \mathcal{G}$ If $\zeta = 1$, since μ^* is a ε stable distribution we have

$$\begin{aligned}
 1871 & |\mu^*(s_v) - \mu^*(s_{\text{base}})P(s_v|s_{\text{base}}, \mu^*)| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1872 & \\
 1873 & \left| \mu^*(s_v) - \mu^*(s_{\text{base}}) \frac{\theta}{\max\{B, \mu^*(s_{\text{base}})\}} \right| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1874 & \\
 1875 & |\mu^*(s_v) - \theta| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1876 & \\
 1877 & \left| \frac{\mu^*(s_v)}{\theta} - 1 \right| \leq \frac{\varepsilon}{\theta|\mathcal{S}|} \leq \frac{\varepsilon}{\theta V} \leq 8\varepsilon, \\
 1878 & \\
 1879 &
 \end{aligned}$$

1880 where we used the fact that $\frac{\theta}{\max\{B, \mu^*(s_{\text{base}})\}} = \mu^*(s_{\text{base}})$. and it follows by definition that $|u_1\left(\frac{\mu^*(s_v)}{\theta}\right) - 1| \leq 8\varepsilon$, since the
 1881 map u_1 is 1-Lipschitz and therefore can only decrease the absolute value on the left. Likewise, if $\zeta = 0$,

$$\begin{aligned}
 1882 & \\
 1883 & \left| \mu^*(s_v) - \sum_{s \in \mathcal{S}} \mu^*(s)P(s_v|s, \mu^*) \right| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1884 & \\
 1885 & |\mu^*(s_v)| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1886 & \\
 1887 & \left| \frac{\mu^*(s_v)}{\theta} \right| \leq \frac{\varepsilon}{\theta|\mathcal{S}|} \leq 8\varepsilon \\
 1888 & \\
 1889 &
 \end{aligned}$$

1890 and once again $u_1\left(\frac{\mu^*(s_v)}{\theta}\right) \leq 8\varepsilon$.

1892 **Case 2: Weighted addition error.** Assume that $G_{\times,+}(\alpha, \beta|v_1, v_2|v) \in \mathcal{G}$, and set $\square := u_\theta(\alpha u_\theta(\mu(v_1)) + \beta u_\theta(\mu(v_2)))$.
 1893 Using the fact that $\|\mu^* - \Gamma_P(\mu^*)\| \leq \frac{\varepsilon}{|\mathcal{S}|}$,

$$\begin{aligned}
 1894 & \\
 1895 & \left| \mu^*(s_v) - \sum_{s \in \mathcal{S}} \mu^*(s)P(s_v|s, \mu^*) \right| \leq \frac{\varepsilon}{|\mathcal{S}|}, \\
 1896 & \\
 1897 & \left| \mu^*(s_v) - \mu^*(s_{\text{base}}) \frac{u_\theta(\alpha u_\theta(\mu(v_1)) + \beta u_\theta(\mu(v_2)))}{\max\{B, \mu^*(s_{\text{base}})\}} \right| \leq \frac{\varepsilon}{|\mathcal{S}|}, \\
 1898 & \\
 1899 & \left| \frac{\mu^*(s_v)}{\theta} - \frac{\square}{\theta} \right| \leq \frac{\varepsilon}{|\mathcal{S}|\theta}, \\
 1900 & \\
 1901 &
 \end{aligned}$$

1902 which implies

$$\left| u_1\left(\frac{\mu^*(s_v)}{\theta}\right) - u_1\left(\alpha u_1\left(\frac{\mu^*(v_1)}{\theta}\right) + \beta u_1\left(\frac{\mu^*(v_2)}{\theta}\right)\right) \right| \leq 8\varepsilon.$$

1906 **Case 3: Brittle comparison gadget.** Finally, we analyze the more involved case of the comparison gadget. Assume
 1907 $G_{<}(|v_1, v_2|v) \in \mathcal{G}$. The stability conditions for s_v yield:

$$\begin{aligned}
 1908 & \\
 1909 & |\mu^*(s_v) - \mu^*(s_{\text{base}})P(s_v|s_{\text{base}}, \mu^*)| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1910 & \\
 1911 & |\mu^*(s_v) - \theta p_{8\varepsilon}(\theta^{-1}u_\theta(\mu^*(v_1)), \theta^{-1}u_\theta(\mu^*(v_2)))| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1912 &
 \end{aligned}$$

1913 We analyze two cases: $u_1(\theta^{-1}\mu^*(v_1)) \geq u_1(\theta^{-1}\mu^*(v_2)) + 8\varepsilon$ and $u_1(\theta^{-1}\mu^*(v_1)) \leq u_1(\theta^{-1}\mu^*(v_2)) - 8\varepsilon$. In the first
 1914 case, we obtain

$$\theta^{-1}u_\theta(\mu^*(v_1)) \geq \theta^{-1}u_\theta(\mu^*(v_2)) + 8\varepsilon,$$

1917 which implies by the definition of $p_{8\varepsilon}$

$$\begin{aligned}
 1918 & |\mu^*(s_v) - \theta| \leq \frac{\varepsilon}{|\mathcal{S}|} \\
 1919 & \\
 1920 & |u_1(\theta^{-1}\mu^*(s_v)) - 1| \leq \frac{\varepsilon}{|\mathcal{S}|\theta} \\
 1921 & \\
 1922 & u_1(\theta^{-1}\mu^*(s_v)) \geq 1 - \frac{\varepsilon}{|\mathcal{S}|\theta} \geq 1 - 8\varepsilon. \\
 1923 & \\
 1924 &
 \end{aligned}$$

1925 In the second case $u_1(\theta^{-1}\mu^*(v_1)) \leq u_1(\theta^{-1}\mu^*(v_2)) - 8\varepsilon$, it follows by a similar analysis that

$$1926 \quad u_1(\theta^{-1}\mu^*(s_v)) \leq \frac{\varepsilon}{|S|\theta} \leq 8\varepsilon.$$

1927
1928
1929 Hence, in the above, we reduced the 8ε -GCIRCUIT problem to the ε -STATDIST problem, completing the proof that
1930 ε -STATDIST is PPAD-hard. The fact that ε -STATDIST is in PPAD on the other hand easily follows from the fact that
1931 ε -STATDIST is the fixed point problem for the (simple) operator Γ_P , reducing it to the END-OF-THE-LINE problem by a
1932 standard construction (Daskalakis et al., 2009).
1933

1934 B.3. Proof of Intractability of FH-MFG

1936 As in the previous section, we reduce any ε -GCIRCUIT problem $(\mathcal{G}, \mathcal{V})$ to the problem $(\varepsilon^2, 2)$ -FH-NASH for some simple
1937 reward $R \in \mathcal{R}^{\text{Sim}}$. Once again let $V = |\mathcal{V}|$.

1938 Associated with each $v \in \mathcal{V}$ we define $s_{v,1}, s_{v,0}, s_{v,\text{base}} \in \mathcal{S}$. The initial distribution is defined as

$$1940 \quad \mu_0(s_{v,\text{base}}) = \frac{1}{V}, \forall v \in \mathcal{V},$$

1941 and we define two actions for each state: $\mathcal{A} = \{a_1, a_0\}$. The state transition probability matrix is given by

$$1942 \quad P(s|s_{v,\text{base}}, a) = \begin{cases} 1, & \text{if } a = a_1, s = s_{v,1}, \\ 1, & \text{if } a = a_0, s = s_{v,0}, \\ 0, & \text{otherwise.} \end{cases}$$

$$1943 \quad P(s_{v,\text{base}}|s, a) = 0, \forall v \in \mathcal{V}, s \in \mathcal{S}, a \in \mathcal{A},$$

1944
1945 and an ε satisfying assignment $p : \mathcal{V} \rightarrow [0, 1]$ will be read by $p(v) = \pi_1^*(a_1|s_{v,\text{base}})$ for the optimal policy $\pi^* = \{\pi_h\}_{h=0}^1$.
1946 We will specify population-dependent rewards $R \in \mathcal{R}^{\text{Simple}}$, since R will not depend on the particular action but only the
1947 state and population distribution, we will concisely denote $R(s, a, \mu) = R(s, \mu)$. It will be the case that

$$1948 \quad R(s_{v,\text{base}}, \mu) = 0, \forall v \in \mathcal{V}, \mu \in \Delta_{\mathcal{S}}.$$

1949 We assign $R(s_{v,1}, \mu) = R(s_{v,0}, \mu) = 0, \forall \mu$ for any vertex v of the generalized circuit that is not the output of any gate in \mathcal{G} .

1950 **Binary assignment gadget.** For any binary assignment gate $G_{\leftarrow}(\zeta|v)$, we assign

$$1951 \quad R(s_{v,1}, \mu) = \zeta,$$

$$1952 \quad R(s_{v,0}, \mu) = 1 - \zeta, \forall \mu \in \Delta_{\mathcal{S}}.$$

1953 **Weighted addition gadget.** For any gate $G_{\times,+}(\alpha, \beta|v_1, v_2|v)$,

$$1954 \quad R(s_{v,1}, \mu) = u_1(u_1(\alpha V\mu(s_{v_1,1}) + \beta V\mu(s_{v_2,1})) - V\mu(s_{v,1})),$$

$$1955 \quad R(s_{v,0}, \mu) = u_1(V\mu(s_{v,1}) - u_1(\alpha V\mu(s_{v_1,1}) + \beta V\mu(s_{v_2,1}))),$$

1956 for all $\mu \in \Delta_{\mathcal{S}}$.

1957 **Brittle comparison gadget.** For any gate $G_{<}(|v_1, v_2|v)$, we define the rewards for states $s_{v,1}, s_{v,0}$ as

$$1958 \quad R(s_{v,1}, \mu) = u_1(V\mu(s_{v_2,1}) - V\mu(s_{v_1,1})),$$

$$1959 \quad R(s_{v,0}, \mu) = u_1(V\mu(s_{v_1,1}) - V\mu(s_{v_2,1})), \forall \mu \in \Delta_{\mathcal{S}}.$$

1960 Now assume that $\pi^* = \{\pi_h^*\}_{h=0}^1$ is a solution to the $(\varepsilon^2, 2)$ -FH-NASH problem and $\mu^* = \Lambda_{P,\mu_0}^2(\pi^*)$, that is, assume that
1961 for all $\pi \in \Pi^2$,

$$1962 \quad V_{P,R}^H(\mu^*, \pi) - V_{P,R}^H(\mu^*, \pi^*) \leq \frac{\varepsilon^2}{V}.$$

1980 Firstly, if μ_1^* is induced by π^* , it holds that $\forall v \in \mathcal{V}$,

$$1981 \mu_1^*(s_{v,\text{base}}) = 0, \quad \mu_1^*(s_{v,1}) = \frac{1}{V} \pi_0^*(s_{v,1}|s_{v,\text{base}}),$$

$$1982 \mu_1^*(s_{v,0}) = \frac{1 - \pi_0^*(s_{v,1}|s_{v,\text{base}})}{V}.$$

1986 Furthermore, a policy $\pi^{\text{br}} \in \Pi_2$ that is the best response to $\mu^* := \{\mu_0^*, \mu_1^*\}$ can be always formulated as:

$$1987 \pi_0^{\text{br}}(a_1|s_{v,\text{base}}) = \begin{cases} 1, & \text{if } R(s_{v,1}, \mu_1^*) > R(s_{v,1}, \mu_1^*), \\ 0, & \text{otherwise} \end{cases}$$

$$1988 \pi_0^{\text{br}}(a_0|s_{v,\text{base}}) = 1 - \pi_0^{\text{br}}(a_1|s_{v,\text{base}}),$$

$$1989 \pi_1^{\text{br}}(a_1|s_{v,\text{base}}) = 1,$$

$$1990 \pi_1^{\text{br}}(a_0|s_{v,\text{base}}) = 0.$$

1995 By the optimality conditions, we will have

$$1996 V_{P,R}^H(\mu^*, \pi^{\text{br}}) - V_{P,R}^H(\mu^*, \pi^*) \leq \frac{\varepsilon^2}{V}.$$

1999 Furthermore, for any $v \in \mathcal{V}$ it holds that

$$2000 V_{P,R}^H(\mu^*, \pi^{\text{br}}) - V_{P,R}^H(\mu^*, \pi^*)$$

$$2001 = \sum_{v \in \mathcal{V}} \mu_0(s_{v,\text{base}}) \left[\max_{s \in \{s_{v,1}, s_{v,0}\}} R(s, \mu_1^*) \right.$$

$$2002 \quad \left. - \pi_0^*(a_1|s_{v,\text{base}}) R(s_{v,1}, \mu_1^*) - \pi_0^*(a_0|s_{v,\text{base}}) R(s_{v,0}, \mu_1^*) \right]$$

$$2003 \geq \frac{1}{V} \max_{s \in \{s_{v,1}, s_{v,0}\}} R(s, \mu_1^*)$$

$$2004 \quad - \frac{1}{V} \pi_0^*(a_1|s_{v,\text{base}}) R(s_{v,1}, \mu_1^*) - \frac{1}{V} \pi_0^*(a_0|s_{v,\text{base}}) R(s_{v,0}, \mu_1^*)$$

2010 as the summands are all positive. We prove that all gate conditions are satisfied case by base. Without loss of generality, we
2011 assume $\varepsilon < 1$ below.

2012 **Case 1.** It follows that for any $v \in \mathcal{V}$ such that $G_{\leftarrow}(\zeta|v) \in \mathcal{G}$, we have

$$2013 \frac{1}{V} - \frac{1}{V} \pi_0^*(a_1|s_{v,\text{base}}) \zeta - \frac{1}{V} \pi_0^*(a_0|s_{v,\text{base}}) (1 - \zeta) \leq \frac{\varepsilon^2}{V}$$

$$2014 1 - \pi_0^*(a_1|s_{v,\text{base}}) \zeta - (1 - \pi_0^*(a_1|s_{v,\text{base}})) (1 - \zeta) \leq \varepsilon^2$$

$$2015 \zeta(1 - 2\pi_0^*(a_1|s_{v,\text{base}})) + \pi_0^*(a_1|s_{v,\text{base}}) \leq \varepsilon^2 \leq \varepsilon.$$

2019 The above implies $\pi_0^*(a_1|s_{v,\text{base}}) \geq 1 - \varepsilon$ if $\zeta = 1$, and if $\zeta = 0$, it implies $\pi_0^*(a_1|s_{v,\text{base}}) \leq \varepsilon$.

2020 **Case 2.** For any $v \in \mathcal{V}$ such that $G_{\times,+}(\alpha, \beta|v_1, v_2|v) \in \mathcal{G}$, denoting in short

$$2021 \square := u_1(\alpha V \mu_1^*(s_{v_1,1}) + \beta V \mu_1^*(s_{v_2,1}))$$

$$2022 = u_1(\alpha \pi_0^*(a_1|s_{v_1,1}) + \beta \pi_0^*(a_1|s_{v_2,1})),$$

$$2023 p_1 := \pi_0^*(a_1|s_{v,\text{base}})$$

$$2024 p_0 := \pi_0^*(a_0|s_{v,\text{base}})$$

2027 we have

$$2028 \frac{1}{V} \max \{ u_1(V \mu_1^*(s_{v,1}) - \square), u_1(\square - V \mu_1^*(s_{v,1})) \}$$

$$2029 \quad - \frac{1}{V} \pi_0^*(a_1|s_{v,\text{base}}) u_1(\square - V \mu_1^*(s_{v,1}))$$

$$2030 \quad - \frac{1}{V} \pi_0^*(a_0|s_{v,\text{base}}) u_1(V \mu_1^*(s_{v,1}) - \square) \leq \varepsilon^2,$$

2035 or equivalently

$$2036 \max \{u_1(p_1 - \square), u_1(\square - p_1)\} - p_1 u_1(\square - p_1) - p_0 u_1(p_1 - \square) \leq \varepsilon^2.$$

2037
2038 First, assume it holds that $p_1 \leq \square$, then:

$$2039 \begin{aligned} 2040 u_1(\square - p_1) - p_1 u_1(\square - p_1) &\leq \varepsilon^2 \\ 2041 (1 - p_1)(\square - p_1) &\leq \varepsilon^2. \end{aligned}$$

2042 The above implies that either $p_1 \geq 1 - \varepsilon$ or $u_1(\square - p_1) \leq \varepsilon$, both cases implying $|\square - p_1| \leq \varepsilon$ since we assume $\square \geq p_1$.
2043 To conclude case 2, assume that $\square < p_1$, then

$$2044 \begin{aligned} 2045 u_1(p_1 - \square) - (1 - p_1)u_1(p_1 - \square) &\leq \varepsilon^2, \\ 2046 p_1(p_1 - \square) &\leq \varepsilon^2, \end{aligned}$$

2047 then either $p_1 \leq \varepsilon$ or $p_1 - \square \leq \varepsilon$, either case implying once again $|\square - p_1| \leq \varepsilon$.

2048 **Case 3.** Finally, for any $v \in \mathcal{V}$ such that $G_{<}(|v_1, v_2|v) \in \mathcal{G}$,

$$2049 \begin{aligned} 2050 \frac{1}{V} \max \{ &u_1(\mu(s_{v_2,1}) - \mu(s_{v_1,1})), u_1(\mu(s_{v_1,1}) - \mu(s_{v_2,1})) \} \\ 2051 &- \frac{1}{V} \pi_0^*(a_1|s_{v,\text{base}})u_1(\mu(s_{v_1,1}) - \mu(s_{v_2,1})) \\ 2052 &- \frac{1}{V} \pi_0^*(a_0|s_{v,\text{base}})u_1(\mu(s_{v_2,1}) - \mu(s_{v_1,1})) \leq \varepsilon \end{aligned}$$

2053 hence once again using the shorthand notation:

$$2054 \begin{aligned} 2055 \Delta &:= V\mu_1^*(s_{v_2,1}) - V\mu_1^*(s_{v_1,1}) = \pi_0^*(a_1|s_{v_2,1}) - \pi_0^*(a_1|s_{v_1,1}) \\ 2056 p_1 &:= \pi_0^*(a_1|s_{v,\text{base}}) \\ 2057 p_0 &:= \pi_0^*(a_0|s_{v,\text{base}}) \end{aligned}$$

2058 we have the inequality:

$$2059 \begin{aligned} 2060 u_1(|\Delta|) - p_1 u_1(\Delta) - p_0 u_1(-\Delta) &\leq \varepsilon^2 \\ 2061 u_1(|\Delta|) - p_1 u_1(\Delta) - (1 - p_1)u_1(-\Delta) &\leq \varepsilon^2. \end{aligned}$$

2062 First assume $\Delta \geq \varepsilon$, then

$$2063 u_1(\Delta)(1 - p_1) \leq \varepsilon^2 \implies 1 - \varepsilon \leq p_1,$$

2064 and conversely if $\Delta \leq -\varepsilon$,

$$2065 u_1(-\Delta)p_1 \leq \varepsilon^2 \implies p_1 \leq \varepsilon,$$

2066 concluding that the comparison gate conditions are ε satisfied for the assignment $v \rightarrow \pi_0^{\text{br}}(a_1|s_{v,\text{base}})$.

2067 The three cases above conclude that $v \rightarrow \pi_0^{\text{br}}(a_1|s_{v,\text{base}})$ is an ε -satisfying assignment for the generalized circuit $(\mathcal{V}, \mathcal{G})$,
2068 concluding the proof that $(\varepsilon_0, 2)$ -FH-NASH is PPAD-hard for some $\varepsilon_0 > 0$. The fact that $(\varepsilon_0, 2)$ -FH-NASH is in PPAD
2069 follows from the fact that the NE is a fixed point of a simple map on space Π_2 , see for instance (Huang et al., 2023).

2070 B.4. Proof of Intractability of 2-FH-LINEAR

2071 Our reduction will be similar to the previous section, however, instead of reducing a ε -GCIRCUIT to an MFG, we will
2072 reduce a 2 player general sum normal form game, 2-NASH, to a finite horizon mean field game with linear rewards with
2073 horizon $H = 2$ (2-FH-LINEAR). Let $\varepsilon > 0, K_1, K_2 \in \mathbb{N}_{>0}, A, B \in \mathbb{R}^{K_1, K_2}$ be given for a 2-NASH problem. We assume
2074 without loss of generality that $K_1 > 1$, as otherwise, the solution of 2-NASH is trivial.

2090 This time, we define finite horizon game with $K_1 + K_2 + 2$ states, denoted $\mathcal{S} := \{s_{\text{base}}^1, s_{\text{base}}^2, s_1^1, \dots, s_{K_1}^1, s_1^2, \dots, s_{K_2}^2\}$.
 2091 Without loss of generality, we can assume $K_1 \leq K_2$. The action set will be defined by $\mathcal{A} = [K_2] = \{1, \dots, K_2\}$. The
 2092 initial state distribution will be given by $\mu_0(s_{\text{base}}^1) = \mu_0(s_{\text{base}}^2) = 1/2$, with $\mu_0(s) = 0$ for all other states. We define the
 2093 transitions for any $s \in \mathcal{S}, a, a' \in \mathcal{A}$ as:

$$\begin{aligned}
 P(s|s_{\text{base}}^1, a) &= \begin{cases} 1, & \text{if } s = s_a^1 \text{ and } a \leq K_1, \\ 1, & \text{if } s = s_a^1 \text{ and } a > K_1, \\ 0, & \text{otherwise.} \end{cases} \\
 P(s|s_{\text{base}}^2, a) &= \begin{cases} 1, & \text{if } s = s_a^2, \\ 0, & \text{otherwise.} \end{cases} \\
 P(s|s_a^1, a') &= \begin{cases} 1, & \text{if } s = s_a^1, \\ 0, & \text{otherwise.} \end{cases} \quad P(s|s_a^2, a') = \begin{cases} 1, & \text{if } s = s_a^2, \\ 0, & \text{otherwise.} \end{cases}
 \end{aligned}$$

2106 Finally, we will define the linear reward function as for all $a \in [K_2]$:

$$\begin{aligned}
 R(s_{\text{base}}^1, a, \mu) &= 0, \\
 R(s_{\text{base}}^2, a, \mu) &= 0, \\
 R(s_a^1, a, \mu) &= \begin{cases} 0, & \text{if } a > K_1, \\ \frac{1}{2} + \frac{1}{2} \sum_{a' \in [K_2]} \mu(s_{a'}^2) A_{a,a'} \end{cases} \\
 R(s_a^2, a, \mu) &= \frac{1}{2} + \frac{1}{2} \sum_{a' \in [K_1]} \mu(s_{a'}^1) B_{a',a}.
 \end{aligned}$$

2118 In words, the states $s_{\text{base}}^1, s_{\text{base}}^2$ represent the two players of the 2-NASH, and an agent starting from one of the initial base
 2119 states $s_{\text{base}}^1, s_{\text{base}}^2$ of the FH-MFG at round $h = 0$ will be placed at $h = 1$ at a state representing the (pure) strategies of each
 2120 player respectively.

2122 Given the game description above, assume $\pi^* = \{\pi_h^*\}_{h=0}^1$ is an ε solution of the 2-FH-LINEAR. Then, it holds for the
 2123 induced distribution $\mu^* := \{\mu_h^*\}_{h=0}^1 = \Lambda_P^H$ that:

$$\begin{aligned}
 \mu_0^* &= \mu_0, \\
 \mu_1^*(s) &= \sum_{s', a' \in \mathcal{S} \times \mathcal{A}} \mu_0(s') \pi^*(a'|s') P(s|s', a') \\
 &= \begin{cases} \frac{1}{2} \pi_0(i|s_{\text{base}}^1), & \text{if } s = s_i^1, \text{ for some } i \in [K_1], \\ \frac{1}{2} \pi_0(i|s_{\text{base}}^2), & \text{if } s = s_i^2, \text{ for some } i \in [K_2], \\ \frac{1}{2} - \frac{1}{2} \sum_{i \in [K_1]} \pi_0(i|s_{\text{base}}^1), & \text{if } s = s_{\text{base}}^1, \\ 0, & \text{otherwise.} \end{cases}
 \end{aligned}$$

2135 By definition of the ε finite horizon Nash equilibrium,

$$\mathcal{E}_{P,R}^H(\pi^*) := \max_{\pi' \in \Pi^H} V_{P,R}^H(\Lambda_P^H(\pi^*), \pi') - V_{P,R}^H(\Lambda_P^H(\pi^*), \pi) \leq \varepsilon,$$

2140 in particular, it holds for any $\pi \in \Pi_2$ that

$$V_{P,R}^H(\mu^*, \pi) - V_{P,R}^H(\mu^*, \pi^*) \leq \varepsilon. \tag{12}$$

2145 By direct computation, the value functions $V_{P,R}^H$ can be written directly in this case for any π :

$$\begin{aligned}
 2146 & \\
 2147 & \\
 2148 & V_{P,R}^H(\boldsymbol{\mu}^*, \boldsymbol{\pi}) = \frac{1}{2} \sum_{a \in [K_1]} \pi_0(a|s_{\text{base}}^1) \left(\frac{1}{2} + \frac{1}{2} \sum_{a' \in [K_2]} \mu_1^*(s_{a'}^2) A_{a,a'} \right) \\
 2149 & \\
 2150 & \\
 2151 & \quad + \frac{1}{2} \sum_{a' \in [K_2]} \pi_0(a'|s_{\text{base}}^2) \left(\frac{1}{2} + \frac{1}{2} \sum_{a \in [K_1]} \mu_1^*(s_a^1) B_{a,a'} \right) \\
 2152 & \\
 2153 & \\
 2154 & = \frac{1}{4} \left(1 + \sum_{a \in [K_1]} \pi_0(a|s_{\text{base}}^1) \right) \\
 2155 & \\
 2156 & \quad + \frac{1}{8} \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \\
 2157 & \\
 2158 & \quad + \frac{1}{8} \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0(a'|s_{\text{base}}^2) \pi_0^*(a|s_{\text{base}}^1) B_{a,a'} \\
 2159 & \\
 2160 & \\
 2161 & \\
 2162 &
 \end{aligned}$$

2163 We analyze two different cases, accounting for a possible imbalance between the strategy spaces of the two players, $[K_1]$
2164 and $[K_2]$.

2165 **Case 1.** Assume $K_1 = K_2$. Then, $V_{P,R}^H(\boldsymbol{\mu}^*, \boldsymbol{\pi})$ simplifies to

$$\begin{aligned}
 2166 & \\
 2167 & V_{P,R}^H(\boldsymbol{\mu}^*, \boldsymbol{\pi}) = \frac{1}{2} + \frac{1}{8} \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \\
 2168 & \\
 2169 & \quad + \frac{1}{8} \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0(a'|s_{\text{base}}^2) \pi_0^*(a|s_{\text{base}}^1) B_{a,a'}. \tag{13} \\
 2170 & \\
 2171 & \\
 2172 &
 \end{aligned}$$

2173 Take an arbitrary mixed strategy $\sigma_1 \in \Delta_{[K_1]}$ and define the policy $\boldsymbol{\pi}_A = \{\pi_{A,h}\}_{h=0} \in \Pi^2$ so that

$$2174 \quad \pi_{A,0}(s_{\text{base}}^1) = \sigma_1, \quad \pi_{A,0}(s_{\text{base}}^2) = \pi_0^*(s_{\text{base}}^2), \quad \pi_{A,1} = \pi_1^*.$$

2176 Then, placing $\boldsymbol{\pi}_A$ in equations (13) and (12), it follows that

$$\begin{aligned}
 2177 & \\
 2178 & \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \sigma_1(a) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \\
 2179 & \\
 2180 & \quad - \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0^*(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \leq 8\varepsilon. \tag{14} \\
 2181 & \\
 2182 &
 \end{aligned}$$

2183 Similarly, for any $\sigma_2 \in \Delta_{[K_2]}$, replacing $\boldsymbol{\pi}$ in equations (13) and (12) with a policy $\boldsymbol{\pi}_B$ such that

$$2184 \quad \pi_{B,0}(s_{\text{base}}^1) = \pi_0^*(s_{\text{base}}^1), \quad \pi_{B,0}(s_{\text{base}}^2) = \sigma_2, \quad \pi_{B,1} = \pi_1^*,$$

2186 we obtain

$$\begin{aligned}
 2187 & \\
 2188 & \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \sigma_2(a) \pi_0^*(a'|s_{\text{base}}^1) B_{a,a'} \\
 2189 & \\
 2190 & \quad - \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0^*(a'|s_{\text{base}}^2) \pi_0^*(a|s_{\text{base}}^1) B_{a,a'} \leq 8\varepsilon. \tag{15} \\
 2191 & \\
 2192 &
 \end{aligned}$$

2193 Hence, the resulting equations (14), (15) imply that in this case the strategy profile $(\pi_0^*(s_{\text{base}}^1), \pi_0^*(s_{\text{base}}^2))$ is a 8ε -Nash
2194 equilibrium for the normal form game defined by matrices A, B .

2196 **Case 2.** Next, we analyze the case when $1 < K_1 < K_2$. If $\sum_{a' \in [K_1]} \pi_0^*(a'|s_{\text{base}}^1) = 0$, then the policy

$$2197 \quad \pi'_0(1|s_{\text{base}}^1) = 1, \quad \pi'_0(s_{\text{base}}^2) = \pi_0^*(s_{\text{base}}^2), \quad \pi'_1 = \pi_1^*.$$

2199

2200 yields an exploitability of at least $1/4$, so by taking ε smaller than $1/4$ we can discard this possibility.

2201 Otherwise, we define a policy $\pi_C = \{\pi_{C,h}\}_{h=0}^1 \in \Pi^2$ such that

$$2202 \quad \pi_{C,0}(a|s_{\text{base}}^1) = \begin{cases} \frac{\pi_0^*(a|s_{\text{base}}^1)}{\sum_{a' \in [K_1]} \pi_0^*(a'|s_{\text{base}}^1)}, & \text{if } a \in [K_1], \\ 0, & \text{otherwise.} \end{cases}$$

$$2203 \quad \pi_{C,0}(s_{\text{base}}^2) = \pi_0^*(s_{\text{base}}^2), \quad \pi_{C,1} = \pi_1^*,$$

2204 and replace π in Equation (12) with π_C to obtain:

$$2205 \quad \frac{1}{4} - \frac{1}{4}S$$

$$2206 \quad + \frac{1}{8}(S^{-1} - 1) \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0^*(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \leq \varepsilon$$

2207 where $S := \sum_{a' \in [K_1]} \pi_0^*(a'|s_{\text{base}}^1) < 1$, hence

$$2208 \quad 1 - S = \sum_{a' \in [K_2] - [K_1]} \pi_0^*(a'|s_{\text{base}}^1) \leq 4\varepsilon.$$

2209 Now for some $\sigma_1 \in \Delta_{[K_1]}$, once again take the policy π_A defined in Case 1, and use Inequality (12) to obtain:

$$2210 \quad \frac{1}{4}(1 - S) + \frac{1}{8} \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \sigma_1(a) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'}$$

$$2211 \quad - \frac{1}{8} \sum_{a \in [K_2]} \sum_{a' \in [K_2]} \pi_0^*(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \leq \varepsilon$$

$$2212 \quad \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \sigma_1(a) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'}$$

$$2213 \quad - \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0^*(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) A_{a,a'} \leq 8\varepsilon.$$

2214 Here, using the definition of π_C , as $\pi_{C,0}(a|s_{\text{base}}^1) \geq \pi_0^*(a|s_{\text{base}}^1)$ for $a \in [K_1]$, we obtain:

$$2215 \quad \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \sigma_1(a) \pi_{C,0}(a'|s_{\text{base}}^2) A_{a,a'}$$

$$2216 \quad - \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_{C,0}(a|s_{\text{base}}^1) \pi_{C,0}(a'|s_{\text{base}}^2) A_{a,a'} \leq 8\varepsilon.$$

2217 Next take π_B as defined above in Case 1 for any arbitrary $\sigma_2 \in \Delta_{[K_2]}$ and use the Inequality 12:

$$2218 \quad \sum_{a' \in [K_2]} \sum_{a \in [K_1]} \sigma_2(a') \pi_0^*(a|s_{\text{base}}^1) B_{a,a'}$$

$$2219 \quad - \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_0^*(a|s_{\text{base}}^1) \pi_0^*(a'|s_{\text{base}}^2) B_{a,a'} \leq 8\varepsilon$$

$$2220 \quad \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \sigma_2(a') \pi_{C,0}(a|s_{\text{base}}^1) B_{a,a'}$$

$$2221 \quad - \sum_{a \in [K_1]} \sum_{a' \in [K_2]} \pi_{C,0}(a|s_{\text{base}}^1) \pi_{C,0}(a'|s_{\text{base}}^2) B_{a,a'} \leq \frac{8\varepsilon}{S} \leq \frac{8\varepsilon}{1 - 4\varepsilon}.$$

2222 Assuming without loss of generality that $\varepsilon < \frac{1}{8}$, it follows that $\pi_{C,0}(s_{\text{base}}^1), \pi_{C,0}(s_{\text{base}}^2)$ is a 16ε solution to the 2-NASH.