

Journal Pre-proof

Graphon Mean-Field Control for Cooperative Multi-Agent Reinforcement Learning

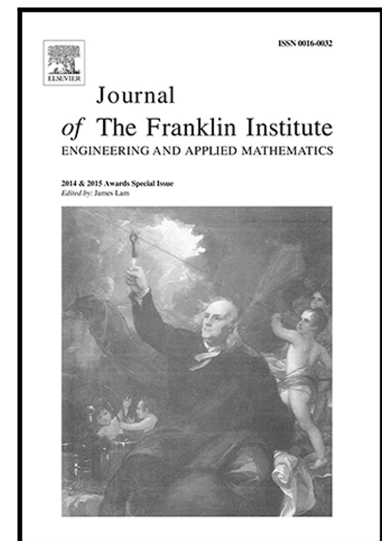
Yuanquan Hu, Xiaoli Wei, Junji Yan, Hengxi Zhang

PII: S0016-0032(23)00548-3
DOI: <https://doi.org/10.1016/j.jfranklin.2023.09.002>
Reference: FI 6391

To appear in: *Journal of the Franklin Institute*

Received date: 8 September 2022
Revised date: 8 May 2023
Accepted date: 1 September 2023

Please cite this article as: Yuanquan Hu, Xiaoli Wei, Junji Yan, Hengxi Zhang, Graphon Mean-Field Control for Cooperative Multi-Agent Reinforcement Learning, *Journal of the Franklin Institute* (2023), doi: <https://doi.org/10.1016/j.jfranklin.2023.09.002>



This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 The Franklin Institute. Published by Elsevier Inc. All rights reserved.

Graphon Mean-Field Control for Cooperative Multi-Agent Reinforcement Learning

Yuanquan Hu, Xiaoli Wei*, Junji Yan, Hengxi Zhang

Abstract

The marriage between mean-field theory and reinforcement learning has shown a great capacity to solve large-scale control problems with homogeneous agents. To break the homogeneity restriction of mean-field theory, a recent interest is to introduce graphon theory to the mean-field paradigm. In this paper, we propose a graphon mean-field control (GMFC) framework to approximate cooperative heterogeneous multi-agent reinforcement learning (MARL) with nonuniform interactions and heterogeneous reward functions and state transition functions among agents and show that the approximate order is of $\mathcal{O}(\frac{1}{\sqrt{N}})$, with N the number of agents. By discretizing the graphon index of GMFC, we further introduce a smaller class of GMFC called block GMFC, which is shown to well approximate cooperative MARL in terms of the value function and the policy. Finally, we design a Proximal Policy Optimization based algorithm for block GMFC that converges to the optimal policy of cooperative MARL. Our empirical studies on several examples demonstrate that our GMFC approach is comparable with the state-of-art MARL algorithms while enjoying better scalability.

Keywords:

Cooperative Multi-Agent Reinforcement Learning, Graphon Theory, Graphon Mean-Field Control, Proximal Policy Optimization

2000 MSC: 60J20, 91A13

1. Introduction

Multi-agent reinforcement learning (MARL) has found various applications in the field of transportation and simulation [50, 1], stock price analysis and trading [32, 31], wireless communication networks [12, 11, 13], and learning behaviors in social dilemmas [33, 28, 34]. MARL, however, becomes intractable due to the complex interactions among agents as the number of agents increases.

A recent tractable approach is a mean-field approach by considering MARL in the regime with a large number of homogeneous agents under weak interactions [20]. According to the

*Corresponding author.

Email addresses: huyq21@mails.tsinghua.edu.cn (Yuanquan Hu), xiaoli.wei@hit.edu.cn (Xiaoli Wei), yan-jj21@mails.tsinghua.edu.cn (Junji Yan), zhanghx20@mails.tsinghua.edu.cn (Hengxi Zhang)

9 number of agents and learning goals, there are three subtle types of mean-field theories for
 10 MARL. The first one is called mean-field MARL (MF-MARL), which refers to the empirical
 11 average of the states or actions of a *finite* population. For example, [52] proposes to approx-
 12 imate interactions within the population of agents by averaging the actions of the overall
 13 population or neighboring agents. [35] proposes a mean-field proximal policy optimization
 14 algorithm for a class of MARL with permutation invariance. The second one is called mean-
 15 field game (MFG), which describes the asymptotic limit of non-cooperative stochastic games
 16 as the number of agents goes to infinity [30, 27, 8]. Recently, a rapidly growing literature
 17 studies MFG for noncooperative MARL either in a model-based way [53, 6, 26] or by a
 18 model-free approach [25, 48, 18, 14, 44]. The third one is called mean-field control (MFC),
 19 which is closely related to MFG yet different from MFG in terms of learning goals. For
 20 cooperative MFC, the Bellman equation for the value function is defined on an enlarged
 21 space of probability measures, and MFC is always reformulated as a new Markov decision
 22 process (MDP) with continuous state-action space. [9] shows the existence of optimal poli-
 23 cies for MFC in the form of mean-field MDP and adapts classical reinforcement learning
 24 (RL) methods to the mean-field setups. [23] approximates MARL by a MFC approach, and
 25 proposes a model-free kernel-based Q-learning algorithm (MFC-K-Q) that enjoys a linear
 26 convergence rate and is independent of the number of agents. [44] presents a model-based
 27 RL algorithm M3-UCRL for MFC with a general regret bound. [2] proposes a unified two-
 28 timescale learning framework for MFG and MFC by tuning the ratio of learning rates of Q
 29 function and the population state distribution. Under the framework of MFC, [41] proposes
 30 locally executable policies such that the resulting discounted sum of average rewards well
 31 approximates the optimal value function over all policies with theoretical guarantee.

32 One restriction of the mean-field theory is that it eliminates the difference among agents
 33 and interactions between agents are assumed to be uniform. However, in many real world
 34 scenarios, strategic interactions between agents are not always uniform and rely on the
 35 relative positions of agents. To develop scalable learning algorithms for multi-agent systems
 36 with heterogeneous agents, one approach is to exploit the local network structure of agents
 37 [45, 37]. Another approach is to consider mean-field systems on large graphs and their
 38 asymptotic limits, which leads to graphon mean-field theory [39]. So far, most existing
 39 works on graphon mean-field theory consider either diffusion processes without learning in
 40 continuous time or non-cooperative graphon mean-field game (GMFG) in discrete time. [3]
 41 considers uncontrolled graphon mean-field systems in continuous time. [17] studies MFG
 42 on an Erdős-Rényi graph. [19] studies the convergence of weighted empirical measures
 43 described by stochastic differential equations. [4] studies propagation of chaos of weakly
 44 interacting particles on general graph sequences. [5] considers general GMFG and studies
 45 ε -Nash equilibria of the multi-agent system by a PDE approach in continuous time. [29]
 46 studies stochastic games on large graphs and their graphon limits. It shows that GMFG
 47 is viewed as a special case of MFG by viewing the label of agents as a component of the
 48 state process. [21, 22] study continuous-time cooperative graphon mean-field systems with
 49 linear dynamics. On the other hand, [7] studies static finite-agent network games and their
 50 associated graphon games. [49] provides a sequential decomposition algorithm to find Nash
 51 equilibria of discrete-time GMFG. [15] constructs a discrete-time learning GMFG framework

52 to analyze approximate Nash equilibria for MARL with nonuniform interactions. However,
 53 little is focused on learning cooperative graphon mean-field systems in discrete time, except
 54 for [42, 43] on particular forms of nonuniform interactions among agents. [43] proves that
 55 when the reward is affine in the state distribution and action distribution, MARL with
 56 nonuniform interactions can still be approximated by classic MFC. [42] considers multi-
 57 class MARL, where agents belonging to the same class are homogeneous. In contrast, we
 58 consider a general discrete-time GMFC framework under which agents are allowed to be
 59 fully heterogeneous and interact non-uniformly on any network captured by a graphon.

60 *Our Work.* In this work, we propose a general discrete-time GMFC framework to approx-
 61 imate cooperative heterogeneous MARL on large graphs by combining classical MFC and
 62 network games. Theoretically, we first show that GMFC can be reformulated as a new
 63 MDP with deterministic dynamics and infinite-dimensional state-action space, hence the
 64 Bellman equation for Q function is established on the space of probability measure ensem-
 65 bles. It shows that GMFC approximates cooperative MARL well in terms of both value
 66 function and optimal policies. The approximation error is at order $\mathcal{O}(1/\sqrt{N})$, where N is
 67 the number of agents. Furthermore, instead of learning infinite-dimensional GMFC directly,
 68 we introduce a smaller class called block GMFC by discretizing the graphon index, which
 69 can be recast as a new MDP with deterministic dynamic and finite-dimensional continuous
 70 state-action space. We show that the optimal policy ensemble learned from block GMFC
 71 is near optimal for cooperative MARL. Using the approach in [38], we develop a proximal
 72 policy optimization (PPO) based algorithm for block GMFC, which, together with approxi-
 73 mation result between block GMFC and cooperative MARL, shows that the proposed PPO
 74 algorithm converges to the optimal policy of MARL with the sample complexity guarantee.
 75 Empirically, our experiments in Section 5 demonstrate that when the number of agents be-
 76 comes large, the mean episode reward of MARL becomes increasingly close to that of block
 77 GMFC, which verifies our theoretical findings. Furthermore, our block GMFC approach
 78 achieves comparable performances with other popular existing MARL algorithms in the
 79 finite-agent setting.

80 *Outline.* The rest of the paper is organized as follows. Section 2 recalls basic notations
 81 of graphons and introduces the setup of cooperative MARL with nonuniform interactions
 82 and its asymptotic limit called GMFC. Section 3 connects cooperative MARL and GMFC,
 83 introduces block GMFC for efficient algorithm design, and builds its connection with coop-
 84 erative MARL. The main theoretical proofs are presented in Section 4. Section 5 tests the
 85 performance of block GMFC experimentally.

86 2. Mean-Field MARL on Dense Graphs

87 2.1. Preliminary: Graphon Theory

88 In the following, we consider a cooperative multi-agent system and its associated mean-
 89 field limit. In this system, each agent is affected by all others, with different agents exerting
 90 different effects on her. This multi-agent system with N agents can be described by a
 91 weighted graph $G_N = (\mathcal{V}_N, \mathcal{E}_N)$, where the vertex set $\mathcal{V}_N = \{1, \dots, N\}$ and the edge set \mathcal{E}_N

92 represent agents and the interactions between agents, respectively. The adjacency matrix
 93 of G_N is represented as $\{\xi_{i,j}^N\}_{1 \leq i,j \leq N}$. To study the limit of the multi-agent system as N
 94 goes to infinity, we adopt the graphon theory introduced in [39] used to characterize the
 95 limit behavior of dense graph sequences. Therefore, throughout the paper, we assume the
 96 graph G_N is dense and leave sparse graphs for future study.

97 In general, a graphon is represented by a bounded symmetric measurable function $W : \mathcal{I} \times \mathcal{I} \rightarrow \mathcal{I}$, with $\mathcal{I} = [0, 1]$. We denote by \mathcal{W} the space of all graphons and equip the space
 98 \mathcal{W} with the cut norm $\|\cdot\|_{\square}$
 99

$$\|W\|_{\square} = \sup_{S,T \subset \mathcal{I}} \left| \int_{S \times T} W(\alpha, \beta) d\alpha d\beta \right|.$$

100 For each weighted graph $G_N = (\mathcal{V}_N, \mathcal{E}_N)$, we consider the correspondence between the
 101 adjacency matrix $\{\xi_{i,j}^N\}$ and a function on $\mathcal{I} \times \mathcal{I}$ with constant value $\xi_{i,j}^N$ on each block
 102 $(\frac{i-1}{N}, \frac{i}{N}] \times (\frac{j-1}{N}, \frac{j}{N}]$. We make the following condition on the strength of interaction $\xi_{i,j}^N$
 103 between agents i and j and the associated W_N .

104 **Condition on W_N and $\xi_{i,j}^N$**

105 1) W_N is a step graphon, that is, $0 \leq W_N \leq 1$ and W_N is a constant on each block
 106 $(\frac{i-1}{N}, \frac{i}{N}] \times (\frac{j-1}{N}, \frac{j}{N}]$:

$$W_N(\alpha, \beta) = W_N\left(\frac{i}{N}, \frac{j}{N}\right), \text{ if } \alpha \in \left(\frac{i-1}{N}, \frac{i}{N}\right], \beta \in \left(\frac{j-1}{N}, \frac{j}{N}\right]. \quad (2.1)$$

2) $\xi_{i,j}^N$ is taken as either

$$\xi_{i,j}^N = W_N\left(\frac{i}{N}, \frac{j}{N}\right) \quad (C1)$$

or

$$\xi_{i,j}^N \sim \text{Bernoulli}\left(W_N\left(\frac{i}{N}, \frac{j}{N}\right)\right). \quad (C2)$$

107 We further assume that the sequence of W_N converges to a graphon W in cut norm as
 108 the number of agents N goes to infinity, which is crucial for the convergence analysis of
 109 cooperative MARL in Section 3.

110 **Assumption 2.1** *The sequence $(W_N)_{N \in \mathbb{N}}$ converges in cut norm to some graphon $W \in \mathcal{W}$*
 111 *such that*

$$\|W_N - W\|_{\square} \rightarrow 0.$$

112 Some common examples of graphons include

113 1) Erdős Rényi: $W(\alpha, \beta) = p$, $0 \leq p \leq 1$, $\alpha, \beta \in \mathcal{I}$;

114 2) Stochastic block model:

$$W(\alpha, \beta) = \begin{cases} p & \text{if } 0 \leq \alpha, \beta \leq 0.5 \text{ or } 0.5 \leq \alpha, \beta \leq 1, \\ q & \text{otherwise,} \end{cases}$$

115 where p represents the intra-community interaction and q the inter-community inter-
116 action;

117 3) Random geometric graphon: $W(\alpha, \beta) = f(\min(|\beta - \alpha|, 1 - |\beta - \alpha|))$, where $f : [0, 0.5] \rightarrow$
118 $[0, 1]$ is a non-increasing function.

119 2.2. Cooperative Heterogeneous MARL

120 In this section, we facilitate the analysis of MARL by considering a particular class of
121 MARL with nonuniform interactions, where each agent interacts with all other agents via
122 the aggregated weighted mean-field effect of the population of all agents.

123 Recall that we use the weighted graph $G_N = (\mathcal{V}_N, \mathcal{E}_N)$ to represent the multi-agent
124 system, in which agents are cooperative and coordinated by a central controller. They
125 share a finite state space \mathcal{S} and take actions from a finite action space \mathcal{A} . We denote by
126 $\mathcal{P}(\mathcal{S})$ and $\mathcal{P}(\mathcal{A})$ the space of all probability measures on \mathcal{S} and \mathcal{A} , respectively. Furthermore,
127 denote by $\mathcal{B}(\mathcal{S})$ the space of all Borel measures on \mathcal{S} .

128 For each agent i , the *neighborhood empirical measure* is given by

$$\mu_t^{i, W_N}(\cdot) := \frac{1}{N} \sum_{j \in \mathcal{V}_N} \xi_{i,j}^N \delta_{s_t^j}(\cdot), \quad (2.2)$$

129 where $\delta_{s_t^j}$ denotes Dirac measure at s_t^j , and (See [15] for more details).

130 At each step $t = 0, 1, \dots$, if agent i , $i \in [N]$ at state $s_t^i \in \mathcal{S}$ takes an action $a_t^i \in \mathcal{A}$, then
131 she will receive a reward

$$r^i(s_t^i, \mu_t^{i, W_N}, a_t^i), \quad i \in [N], \quad (2.3)$$

132 where $r^i : \mathcal{S} \times \mathcal{B}(\mathcal{S}) \times \mathcal{A} \rightarrow \mathbb{R}$, $i \in [N]$, and she will change to a new state s_{t+1}^i according to
133 a transition probability such that

$$s_{t+1}^i \sim P^i(\cdot | s_t^i, \mu_t^{i, W_N}, a_t^i), \quad i \in [N], \quad s_0^i \sim \mu \in \mathcal{P}(\mathcal{S}), \quad (2.4)$$

134 where $P^i : \mathcal{S} \times \mathcal{B}(\mathcal{S}) \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$, $i \in [N]$.

135 (2.3)-(2.4) indicate that the reward and the transition probability of agent i at time
136 t depend on both her individual information (s_t^i, a_t^i) and neighborhood empirical measure
137 μ_t^{i, W_N} .

138 Furthermore, the policy is assumed to be stationary for simplicity and takes the Marko-
139 vian form

$$a_t^i \sim \pi^i(\cdot | s_t^i, \mu_t^{i, W_N}) \in \mathcal{P}(\mathcal{A}), \quad i \in [N], \quad (2.5)$$

140 which maps agent i 's state to a randomized action. (2.5) is called global policy since the
141 policy of agent i depends on both her own state and the aggregate information of the whole
142 population. For each agent i , the space of all global policies is denoted as Π .

143 **Remark 2.2** *It is computationally expensive to collect the aggregate information of the*
 144 *whole population in many practical scenarios. Considering the costly collection of the ag-*
 145 *gregation information of the whole population, one can restrict the policy to be in a local*
 146 *manner, that is, the policy that the agent i can execute depends solely on her own state*
 147 *information:*

$$a_t^i \sim \pi^i(\cdot | s_t^i) \in \mathcal{P}(\mathcal{A}), \quad i \in [N].$$

148 *This has been studied in [41] for standard MFC. Precisely, [41] designs locally executable*
 149 *policies such that the resulting discounted sum of average rewards well approximates the*
 150 *optimal value function over all policies. We expect that a similar result holds for GMFC.*

151 **Remark 2.3** *When $\xi_{ij}^N \equiv 1$, $r^i \equiv r$, $P^i \equiv P$, $i, j \in [N]$, it corresponds to classical mean-*
 152 *field theory with uniform interactions [9, 23]. Furthermore, our framework is flexible enough*
 153 *to include the nonuniform interactions of actions via $\nu_t^{i, W_N} = \frac{1}{N} \sum_{j \in \mathcal{V}_N} \xi_{ij}^N \delta_{a_t^j}(\cdot)$.*

154 The expected discounted accumulated reward of agent i is

$$J_{N,i}(\mu, \pi^1, \dots, \pi^N) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^i(s_t^i, \mu_t^{i, W_N}, a_t^i) \mid s_0^i \sim \mu, a_t^i \sim \pi^i(\cdot | s_t^i, \mu_t^{i, W_N}) \right], \quad (2.6)$$

155 subject to (2.2)-(2.5) with a discount factor $\gamma \in (0, 1)$.

156 The objective of this cooperative multi-agent system (2.2)-(2.5) is to find Pareto opti-
 157 mality given in the Definition 2.4 below.

158 **Definition 2.4 (Pareto Optimality)** $(\pi^{1,*}, \dots, \pi^{N,*}) \in \Pi^N$ *is called Pareto optimality*
 159 *for the multi-agent system (2.2)-(2.5) if there does not exist $(\pi^1, \dots, \pi^N) \in \Pi^N$ such that*

$$\begin{aligned} \forall 1 \leq i \leq N, & \quad J_{N,i}(\mu, \pi^1, \dots, \pi^N) \geq J_{N,i}(\mu, \pi^{1,*}, \dots, \pi^{N,*}), \\ \exists 1 \leq i \leq N, & \quad J_{N,i}(\mu, \pi^1, \dots, \pi^N) > J_{N,i}(\mu, \pi^{1,*}, \dots, \pi^{N,*}). \end{aligned}$$

160 To study Pareto optimality, we introduce the expected discounted accumulated reward
 161 averaged over all agents, i.e.,

$$\begin{aligned} V_N(\mu) &= \sup_{(\pi^1, \dots, \pi^N) \in \Pi^N} J_N(\mu, \pi^1, \dots, \pi^N) \\ &:= \sup_{(\pi^1, \dots, \pi^N) \in \Pi^N} \frac{1}{N} \sum_{i=1}^N J_{N,i}(\mu, \pi^1, \dots, \pi^N), \end{aligned} \quad (2.7)$$

162 subject to (2.2)-(2.5). Let $(\pi^{1,*}, \dots, \pi^{N,*}) \in \arg \max_{(\pi^1, \dots, \pi^N) \in \Pi^N} J_N(\mu, \pi^1, \dots, \pi^N)$, then $(\pi^{1,*}, \dots, \pi^{N,*})$

163 is shown to be a Pareto optimality in Definition 2.4. Therefore, searching for Pareto opti-
 164 mality of cooperative MARL amounts to solving the optimal policy of (2.7). However, it is
 165 always difficult to exactly obtain the optimal policy of cooperative MARL. We consider a
 166 weak notion of ε -Pareto optimality.

167 **Definition 2.5 (ε -Pareto Optimality)** $(\pi_\varepsilon^{1,*}, \dots, \pi_\varepsilon^{N,*}) \in \Pi^N$ is called ε -Pareto optimal-
 168 ity for the multi-agent system (2.2)-(2.5) if there does not exist $(\pi^1, \dots, \pi^N) \in \Pi^N$ such
 169 that

$$\begin{aligned} \forall 1 \leq i \leq N, J_{N,i}(\mu, \pi^1, \dots, \pi^N) &\geq J_{N,i}(\mu, \pi_\varepsilon^{1,*}, \dots, \pi_\varepsilon^{N,*}) + \varepsilon, \\ \exists 1 \leq i \leq N, J_{N,i}(\mu, \pi^1, \dots, \pi^N) &> J_{N,i}(\mu, \pi_\varepsilon^{1,*}, \dots, \pi_\varepsilon^{N,*}) + \varepsilon. \end{aligned}$$

170 For any $\varepsilon > 0$, let $(\pi_\varepsilon^{1,*}, \dots, \pi_\varepsilon^{N,*}) \in \Pi^N$ such that

$$J_N(\mu, \pi_\varepsilon^{1,*}, \dots, \pi_\varepsilon^{N,*}) \geq \sup_{(\pi^1, \dots, \pi^N) \in \Pi^N} J_N(\mu, \pi^1, \dots, \pi^N) - \varepsilon, \quad (2.8)$$

171 then $(\pi_\varepsilon^{1,*}, \dots, \pi_\varepsilon^{N,*}) \in \Pi^N$ is an ε -Pareto Optimality in Definition 2.5.

172 2.3. Graphon Mean-Field Control

173 We expect the cooperative MARL (2.2)-(2.7) to become a GMFC problem as $N \rightarrow \infty$.
 174 In GMFC, there is a continuum of agents $\alpha \in \mathcal{I}$, and each agent with the index $\alpha \in \mathcal{I}$
 175 follows

$$s_0^\alpha \sim \mu^\alpha, \quad a_t^\alpha \sim \pi^\alpha(\cdot | s_t^\alpha, \mu_t^{\alpha,W}), \quad s_{t+1}^\alpha \sim P^\alpha(\cdot | s_t^\alpha, \mu_t^{\alpha,W}, a_t^\alpha), \quad (2.9)$$

176 where $\mu_t^\alpha = \mathcal{L}(s_t^\alpha)$, $\alpha \in \mathcal{I}$ denotes the probability distribution of s_t^α , and $\mu_t^{\alpha,W}$ is defined as
 177 the *neighborhood mean-field measure* of agent α :

$$\mu_t^{\alpha,W} = \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^\beta d\beta \in \mathcal{B}(\mathcal{S}), \quad (2.10)$$

178 with the graphon W given in Assumption 2.1.

179 To ease the sequel analysis, define the space of state distribution ensembles $\mathcal{M} :=$
 180 $\mathcal{P}(\mathcal{S})^{\mathcal{I}} := \{f : \mathcal{I} \rightarrow \mathcal{P}(\mathcal{S})\}$ and the space of policy ensembles $\mathcal{II} := \mathcal{P}(\mathcal{A})^{\mathcal{S} \times \mathcal{I}}$. Then
 181 $\boldsymbol{\mu} := (\mu^\alpha)_{\alpha \in \mathcal{I}}$ and $\boldsymbol{\pi} := (\pi^\alpha)_{\alpha \in \mathcal{I}}$ are elements in \mathcal{M} and \mathcal{II} , respectively.

182 The objective of GMFC is to maximize the expected discounted accumulated reward
 183 averaged over all agents $\alpha \in \mathcal{I}$

$$\begin{aligned} V(\boldsymbol{\mu}) : &= \sup_{\boldsymbol{\pi} \in \mathcal{II}} J(\boldsymbol{\mu}, \boldsymbol{\pi}) \\ &= \sup_{\boldsymbol{\pi} \in \mathcal{II}} \int_{\mathcal{I}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^\alpha(s_t^\alpha, \mu_t^{\alpha,W}, a_t^\alpha) \mid s_0^\alpha \sim \mu^\alpha, a_t^\alpha \sim \pi^\alpha(\cdot | s_t^\alpha, \mu_t^{\alpha,W}) \right] d\alpha. \end{aligned} \quad (2.11)$$

184 3. Main Results

185 3.1. Reformulation of GMFC

186 In this section, we show that GMFC (2.9)-(2.11) can be reformulated as a MDP with
 187 deterministic dynamics and continuous state-action space $\mathcal{M} \times \mathcal{II}$.

188 **Theorem 3.1** GMFC (2.9)-(2.11) can be reformulated as

$$V(\boldsymbol{\mu}) = \sup_{\boldsymbol{\pi} \in \Pi} \sum_{t=0}^{\infty} \gamma^t R(\boldsymbol{\mu}_t, \boldsymbol{\pi}(\boldsymbol{\mu}_t)), \quad (3.1)$$

189 subject to

$$\mu_{t+1}^{\alpha}(\cdot) = \Phi^{\alpha}(\boldsymbol{\mu}_t, \boldsymbol{\pi}(\boldsymbol{\mu}_t))(\cdot), \quad t \in \mathbb{N}, \quad \mu_0^{\alpha} = \mu^{\alpha}, \quad \alpha \in \mathcal{I}, \quad (3.2)$$

190 where the aggregated reward $R : \mathcal{M} \times \Pi \rightarrow \mathbb{R}$ and the aggregated transition dynamics $\Phi :$
191 $\mathcal{M} \times \Pi \rightarrow \mathcal{M}$ are given by

$$R(\boldsymbol{\mu}, \boldsymbol{\pi}(\boldsymbol{\mu})) = \int_{\mathcal{I}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r^{\alpha}(s, a, \mu^{\alpha, W}) \pi^{\alpha}(a|s, \mu^{\alpha, W}) \mu^{\alpha}(s) d\alpha, \quad (3.3)$$

$$\Phi^{\alpha}(\boldsymbol{\mu}, \boldsymbol{\pi}(\boldsymbol{\mu}))(\cdot) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} P^{\alpha}(\cdot|s, \mu^{\alpha, W}, a) \pi^{\alpha}(a|s, \mu^{\alpha, W}) \mu^{\alpha}(s). \quad (3.4)$$

192 The proof of Theorem 3.1 is similar to the proof of Lemma 2.2 in [24]. So we omit it here.
193 (3.4) and (3.2) indicate the evolution of the state distribution ensemble $\boldsymbol{\mu}_t$ over time.
194 That is, under the fixed policy ensemble $\boldsymbol{\pi}$, the state distribution μ_{t+1}^{α} of agent α at time $t+1$
195 is fully determined by the policy ensemble $\boldsymbol{\pi}$ and the state distribution ensemble $\boldsymbol{\mu}_t$ at time
196 t . Note that the change of *population state distribution ensemble* will affect *neighborhood*
197 *mean-field measure*. In turn, the change of *neighborhood mean-field measure* will have an
198 influence on *population state distribution ensemble*.

199 With the reformulation in Theorem 3.1, the associated Q function starting from $(\boldsymbol{\mu}, \boldsymbol{\pi}) \in$
200 $\mathcal{M} \times \Pi$ is defined as

$$Q(\boldsymbol{\mu}, \boldsymbol{\pi}) = R(\boldsymbol{\mu}, \boldsymbol{\pi}(\boldsymbol{\mu})) + \sup_{\boldsymbol{\pi}' \in \Pi} \left[\sum_{t=1}^{\infty} \gamma^t R(\boldsymbol{\mu}_t, \boldsymbol{\pi}'(\boldsymbol{\mu}_t)) \mid s_0^{\alpha} \sim \mu^{\alpha}, a_0^{\alpha} \sim \pi^{\alpha}(\cdot|s_0^{\alpha}, \mu^{\alpha, W}) \right] \quad (3.5)$$

201 Hence its Bellman equation is given by

$$Q(\boldsymbol{\mu}, \boldsymbol{\pi}) = R(\boldsymbol{\mu}, \boldsymbol{\pi}(\boldsymbol{\mu})) + \gamma \sup_{\boldsymbol{\pi}' \in \Pi} Q(\Phi(\boldsymbol{\mu}, \boldsymbol{\pi}(\boldsymbol{\mu})), \boldsymbol{\pi}'). \quad (3.6)$$

202 **Remark 3.2** (Label-state formulation) GMFC (2.9)-(2.11) can be viewed as a classical MFC
203 with extended state space $\mathcal{S} \times \mathcal{I}$, action space \mathcal{A} , policy $\tilde{\pi} \in \mathcal{P}(\mathcal{A})^{\mathcal{S} \times \mathcal{I}}$, mean-field information
204 $\tilde{\mu} \in \mathcal{P}(\mathcal{S} \times \mathcal{I})$, reward $\tilde{r}((s, \alpha), \tilde{\mu}, a) := r((s, \alpha), \int_{\mathcal{I}} W(\alpha, \beta) \tilde{\mu}(\cdot, \beta) d\beta, a)$, transition dynamics
205 of (\tilde{s}_t, α_t) such that

$$\tilde{s}_{t+1} \sim P(\cdot|(\tilde{s}_t, \alpha_t), \tilde{a}_t, \int_{\mathcal{I}} W(\alpha_t, \beta) \tilde{\mu}_t(\cdot, \beta) d\beta), \quad \alpha_{t+1} = \alpha_t, \quad \tilde{a}_t \sim \tilde{\pi}(\cdot|\tilde{s}_t, \alpha_t, \int_{\mathcal{I}} W(\alpha_t, \beta) \tilde{\mu}_t(\cdot, \beta) d\beta),$$

206 with the initial condition $\tilde{s}_0 \sim \mu_0$, $\tilde{\alpha}_0 \sim Unif(0, 1)$. It is worth pointing out such a label-
207 state formulation has also been studied in GMFG [29, 15].

208 *3.2. Approximation*

209 In this section, we show that GMFC (2.9)-(2.11) provides a good approximation for the
 210 cooperative multi-agent system (2.2)-(2.7) in terms of the value function and the optimal
 211 policy ensemble. To do this, the following assumptions on W , P , r , and $\boldsymbol{\pi}$ are needed.

212 **Assumption 3.3 (graphon W)** *There exists $L_W > 0$ such that for all $\alpha, \alpha', \beta, \beta' \in \mathcal{I}$*

$$|W(\alpha, \beta) - W(\alpha', \beta')| \leq L_W \cdot (|\alpha - \alpha'| + |\beta - \beta'|).$$

213 Assumption 3.3 is common in graphon mean-field theory [21, 15, 29]. Indeed, the Lips-
 214 chitz continuity assumption on W in Assumption 3.3 can be relaxed to piecewise Lipschitz
 215 continuity on W .

216 **Assumption 3.4 (transition probability P)** *There exists $L_P > 0$ and $\tilde{L}_P > 0$ such that*
 217 *for any $\alpha, \beta \in \mathcal{I}$, all $s \in \mathcal{S}$, $a \in \mathcal{A}$, $\mu_1, \mu_2 \in \mathcal{B}(\mathcal{S})$*

$$\|P^\alpha(\cdot|s, \mu_1, a) - P^\beta(\cdot|s, \mu_2, a)\|_1 \leq L_P \cdot \|\mu_1 - \mu_2\|_1 + \tilde{L}_P \cdot |\alpha - \beta|,$$

218 where $\|\cdot\|_1$ denotes L^1 norm here and throughout the paper.

219 **Assumption 3.5 (reward r)** *There exist $M_r > 0$, $L_r > 0$ and $\tilde{L}_r > 0$ such that for all*
 220 *$s \in \mathcal{S}$, $a \in \mathcal{A}$, $\mu_1, \mu_2 \in \mathcal{B}(\mathcal{S})$,*

$$|r^\alpha(s, \mu, a)| \leq M_r, \quad |r^\alpha(s, \mu_1, a) - r^\beta(s, \mu_2, a)| \leq L_r \cdot \|\mu_1 - \mu_2\|_1 + \tilde{L}_r \cdot |\alpha - \beta|.$$

221 **Assumption 3.6 (policy $\boldsymbol{\pi}$)** *There exists $L_\Pi > 0$ and $\tilde{L}_\Pi > 0$ such that for any policy*
 222 *ensemble $\boldsymbol{\pi} := (\pi^\alpha)_{\alpha \in \mathcal{I}} \in \boldsymbol{\Pi}$ is Lipschitz continuous, that is, for any $\alpha, \beta \in \mathcal{I}$ and $\mu_1, \mu_2 \in$*
 223 *$\mathcal{B}(\mathcal{S})$,*

$$\max_{s \in \mathcal{S}} \|\pi^\alpha(\cdot|s, \mu_1) - \pi^\beta(\cdot|s, \mu_2)\|_1 \leq L_\Pi \cdot \|\mu_1 - \mu_2\|_1 + \tilde{L}_\Pi \cdot |\alpha - \beta|.$$

224 Assumptions 3.3-3.6 state that W, P, r and $\boldsymbol{\pi}$ are Lipschitz continuous with respect to
 225 both the index of the agent and the neighborhood mean-field measure. The distance between
 226 indexes $|\alpha - \beta|$ measures the similarity of agents. If P, r and $\boldsymbol{\pi}$ are identical, Assumptions
 227 3.4-3.6 are commonly used to bridge the multi-agent system and classical mean-field theory
 228 [23, 41, 42, 43].

229 To show approximation properties of GMFC in the large-scale multi-agent system, we
 230 need to relate policy ensembles of GMFC to policies of the multi-agent system. On one
 231 hand, one can see that any $\boldsymbol{\pi} \in \boldsymbol{\Pi}$ leads to a N -agent policy tuple $(\pi^1, \dots, \pi^N) \in \Pi^N$ with

$$\Gamma^N : \boldsymbol{\Pi} \ni \boldsymbol{\pi} \mapsto (\pi^1, \dots, \pi^N) \in \Pi^N, \quad \text{with } \pi^i := \pi^{\frac{i}{N}}. \quad (3.7)$$

232 On the other hand, any N -agent policy tuple $(\pi^1, \dots, \pi^N) \in \Pi^N$ can be seen as a step
 233 policy ensemble $\boldsymbol{\pi}^N$ in $\boldsymbol{\Pi}$:

$$\boldsymbol{\pi}^{N, \alpha} := \sum_{i=1}^N \pi^i \mathbf{1}_{\alpha \in (\frac{i-1}{N}, \frac{i}{N}]} \in \boldsymbol{\Pi}. \quad (3.8)$$

234 Similarly, any N -agent reward tuple (r^1, \dots, r^N) can be regarded as a step reward function
 235 of GMFC:

$$r^{N,\alpha} := \sum_{i=1}^N r^i \mathbf{1}_{\alpha \in (\frac{i-1}{N}, \frac{i}{N}]}. \quad (3.9)$$

236

237 **Theorem 3.7 (Approximate Pareto Property)** *Assume Assumptions 2.1, 3.3, 3.4, 3.5*
 238 *and 3.6. Then under either the condition (C1) or (C2), we have for any initial distribution*
 239 $\mu \in \mathcal{P}(\mathcal{S})$

$$|V_N(\mu) - V(\mu)| \rightarrow 0, \text{ as } N \rightarrow \infty. \quad (3.10)$$

240 *Moreover, if the graphon convergence in Assumption 2.1 is at rate $\mathcal{O}(\frac{1}{\sqrt{N}})$, then $|V_N(\mu) -$*
 241 *$V(\mu)| = \mathcal{O}(\frac{1}{\sqrt{N}})$. As a consequence, for any $\varepsilon > 0$, there exists an integer N_ε such that*
 242 *when $N \geq N_\varepsilon$, the optimal policy ensemble of GMFC denoted as π^* (if it exists) provides*
 243 *an ε -Pareto optimality $(\pi^{1,*}, \dots, \pi^{N,*}) := \Gamma^N(\pi^*)$ for the multi-agent system (2.7), with Γ^N*
 244 *defined in (3.7).*

245 Theorem 3.7 implies that if we could compute an algorithm to learn the optimal policy
 246 ensemble of GMFC, then the learned optimal policy ensemble is close to the optimal policy of
 247 MARL. Directly learning the optimal policy of GMFC, however, will lead to high complexity
 248 due to the infinite-dimensional feature of μ and π . Instead, we will introduce a smaller class
 249 of GMFC with a lower dimension in the next section, which enables a scalable algorithm.

250 3.3. Algorithm Design and Convergence Analysis

251 This section will show that discretizing the graphon index $\alpha \in \mathcal{I}$ of GMFC enables to
 252 approximate Q function in (3.6) by an approximated Q function in (3.11) below defined on
 253 a smaller space, which is critical for designing efficient learning algorithms.

254 Precisely, we choose uniform grids $\alpha_m \in \mathcal{I}_M := \{\frac{m}{M}, 0 \leq m \leq M\}$ for simplicity, and
 255 approximate each agent $\alpha \in \mathcal{I}$ by the nearest $\alpha_m \in \mathcal{I}_M$ close to it. Introduce $\tilde{\mathcal{M}}_M :=$
 256 $\mathcal{P}(\mathcal{S})^{\mathcal{I}_M}$, $\tilde{\Pi}_M := \mathcal{P}(\mathcal{A})^{\mathcal{S} \times \mathcal{I}_M}$. Meanwhile, $\tilde{\mu} := (\tilde{\mu}^{\alpha_m})_{m \in [M]} \in \tilde{\mathcal{M}}_M$ and $\tilde{\pi} := (\tilde{\pi}^{\alpha_m})_{m \in [M]} \in$
 257 $\tilde{\Pi}_M$ can be viewed as a piecewise constant state distribution ensemble in \mathcal{M} and a piecewise
 258 constant policy ensemble in Π , respectively. Our arguments can be easily generalized to
 259 nonuniform grids.

260 Consequently, instead of performing algorithms according to (3.6) with a continuum of
 261 graphon labels directly, we work with GMFC with M blocks called **block GMFC**, in which
 262 agents in the same block are homogeneous. The Bellman equation for Q function of block
 263 GMFC is given by

$$\tilde{Q}(\tilde{\mu}, \tilde{\pi}) = \tilde{R}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu})) + \gamma \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \tilde{Q}(\tilde{\Phi}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu})), \tilde{\pi}'), \quad (3.11)$$

264 where the neighborhood mean-field measure, the aggregated reward $\tilde{R} : \tilde{\mathcal{M}}_M \times \tilde{\Pi}_M \rightarrow \mathbb{R}$
 265 and the aggregated transition dynamics $\tilde{\Phi} : \tilde{\mathcal{M}}_M \times \tilde{\Pi}_M \rightarrow \tilde{\mathcal{M}}_M$ are given by

$$\tilde{\mu}^{\alpha_m, W} = \frac{1}{M} \sum_{m'=0}^{M-1} W(\alpha_m, \alpha_{m'}) \tilde{\mu}^{\alpha_{m'}}, m \in [M], \quad (3.12)$$

$$\tilde{R}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu})) = \frac{1}{M} \sum_{m=0}^{M-1} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r^{\alpha_m}(s, a, \tilde{\mu}^{\alpha_m, W}) \tilde{\mu}^{\alpha_m}(s) \tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}^{\alpha_m, W}), \quad (3.13)$$

$$\tilde{\Phi}^{\alpha_m}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu}))(\cdot) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} P^{\alpha_m}(\cdot|s, a, \tilde{\mu}^{\alpha_m, W}) \tilde{\mu}^{\alpha_m}(s) \tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}^{\alpha_m, W}). \quad (3.14)$$

266 We see from (3.11) that block GMFC is a MDP with deterministic dynamics $\tilde{\Phi}$ and
 267 continuous state-action space $\tilde{\mathcal{M}}_M \times \tilde{\Pi}_M$. The following Theorem shows that there exists
 268 an optimal policy ensemble of block GMFC in $\tilde{\Pi}_M$.

269 **Theorem 3.8 (Existence of Optimal Policy Ensemble)** *Given Assumptions 3.4, 3.5,*
 270 *assume $\gamma \cdot (1 + L_P + L_\Pi) < \infty$, then for any fixed integer $M > 0$, there exists an $\tilde{\pi}^* \in \tilde{\Pi}_M$*
 271 *that maximize $\tilde{Q}(\tilde{\mu}, \tilde{\pi})$ in (3.11) for any $\tilde{\mu} \in \tilde{\mathcal{M}}_M$.*

272 Furthermore, we show that with sufficiently fine partitions of the graphon index \mathcal{I} , i.e.,
 273 M is sufficiently large, block GMFC (3.11)-(3.14) well approximates the multi-agent system
 274 in Section 2.2.

275 **Theorem 3.9** *Assume $\gamma \cdot (1 + L_P + L_\Pi) < \infty$ and Assumptions 2.1, 3.3, 3.4, 3.5 and*
 276 *3.6. Under either (C1) or (C2), for any $\varepsilon > 0$, there exists $N_\varepsilon, M_\varepsilon$ such that for $N \geq N_\varepsilon$,*
 277 *the optimal policy ensemble $\tilde{\pi}^*$ of block GMFC (3.11) with M_ε blocks provides an ε -Pareto*
 278 *optimality $(\tilde{\pi}^{1,*}, \dots, \tilde{\pi}^{N,*}) := \Gamma^N(\tilde{\pi}^*)$ for the multi-agent system (2.7) with N agents.*

279 Theorem 3.9 shows that the optimal policy ensemble of block GMFC is near-optimal
 280 for all sufficiently large multi-agent systems, meaning that block GMFC provides a good
 281 approximation for the multi-agent system. Therefore, If we could develop an algorithm for
 282 block GMFC to extract an optimal policy ensemble of block GMFC, then the extracted
 283 policy is near optimal for MARL.

284 When model parameters P^α, r^α and W are known, one can easily extract the optimal
 285 policy based on Bellman equation. If any of these model parameters P^α, r^α and W
 286 unknown, we take a model-free approach. The key issue is to handle population state
 287 distribution ensemble $\tilde{\mu}$, which is an input of \tilde{Q} function in (3.11). We assume that we
 288 have a block GMFC simulator $\mathcal{G}(\tilde{\mu}, \tilde{\pi}) = (\tilde{\mu}', \tilde{R})$. That is, for any pair of population state
 289 distribution ensemble and policy ensemble $(\tilde{\mu}, \tilde{\pi})$, we can sample the aggregated reward \tilde{R}
 290 and the next population state distribution ensemble $\tilde{\mu}'$. To learn the optimal policy of block
 291 GMFC, one can adopt any existing techniques for standard MFC, such as a kernel-based Q
 292 learning method in [23] and a uniform discretization method in [9].

293 **Remark 3.10** *If we can only observe the state of agent $\alpha_m \in \mathcal{I}_M$ and do not have access to*
 294 *population state distribution ensemble, we can estimate $\tilde{\mu}^{\alpha_m}$ following [2] or [42]. However,*

295 different from [2] and [42], we also need to estimate $\tilde{\mu}^{\alpha_m, W}$ due to the graphon structure W
 296 and leave it for future study.

297 We choose to adapt DRL algorithm neural Proximal Policy Optimization (PPO) [47, 38]
 298 to block GMFC given in Algorithm 1. Following Corollary 4.11 in [38], together with
 299 Theorem 3.9, we can state the global convergence of neural PPO for block GMFC. Since
 300 assumptions that make the result hold are similar as [38], we do not state these assumptions
 301 here.

Algorithm 1 Neural PPO for block GMFC

Input Width of neural network M , radius of constraint R , number of SGD and TD iterations T , number of PPO iteration K , penalty parameter β

Initialize

for $k = 0$ to $K - 1$ **do**

set temperature parameter $\tau_{k+1} \leftarrow \frac{\beta\sqrt{K}}{k+1}$ and penalty parameter $\beta_k \leftarrow \beta\sqrt{K}$.

Sample $(\tilde{\mu}_t, \tilde{\pi}_t, \tilde{R}_t, \tilde{\mu}'_t, \tilde{\pi}'_t)_{t=1}^T$ with $\tilde{\pi}_0 \sim \Pi^0(\cdot|\tilde{\mu})$, $\tilde{\mu}'_t = \tilde{\Phi}(\tilde{\mu}_t, \tilde{\pi}_t)$, $\tilde{\pi}_t \sim \Pi^{\theta_k}(\cdot|\tilde{\mu}_t)$.

Solve for Q function parameterized by neural network $Q_{\omega_k} = NN(\omega_k, M)$ using the TD update.

Solve for energy function parameterized by neural network $f_{\theta_{k+1}} = NN(\theta_{k+1}, M)$ using the SGD update.

Update policy: $\Pi^{\theta_k} \propto \exp(\tau_{k+1}^{-1} f_{\theta_{k+1}})$.

end for

Theorem 3.11 Suppose that Assumptions 2.1, 3.3, 3.4, 3.5 and 3.6 hold. Further assume $\gamma \cdot (1 + L_{\Pi} + L_P) < 1$. Furthermore, suppose that the width of neural network is sufficiently large. For any $\varepsilon > 0$, there exists M_{ε} and N_{ε} such that for any $M \geq M_{\varepsilon}$ and $N \geq N_{\varepsilon}$, and the policy attained by Algorithm 1 denoted as π_{PPO}

$$|J_N(\mu; \pi^{1,*}, \dots, \pi^{N,*}) - \tilde{J}^M(\mu; \pi_{PPO})| \leq \frac{C}{\sqrt{K}} + \bar{C}\varepsilon, \quad (3.15)$$

302 where J_N and \tilde{J}^M are given in (2.7) and (4.7) respectively, K is the number of iteration, C
 303 and \bar{C} are constants.

304 By setting $K = \frac{C}{\varepsilon^2}$, the optimal empirical value function of MARL is approximated by the
 305 value function of block GMFC under the learned policy in Algorithm 1 with the error $\mathcal{O}(\varepsilon)$.
 306 In other words, Theorem 3.11 states that, with a sample complexity of $\mathcal{O}(\frac{1}{\varepsilon^2})$, Algorithm 1
 307 generates a $\mathcal{O}(\varepsilon)$ -Pareto optimality of cooperative MARL.

308 To evaluate the performance of Algorithm 1 and to validate our theoretical findings, we
 309 describe the deployment of block GMFC in the multi-agent system in Algorithm 2, which
 310 we call it **N-agent GMFC**.

Algorithm 2 N-agent GMFC

Input Initial state distribution μ_0 , number of agents N , episode length T , the learned policy $\tilde{\pi} \in \tilde{\Pi}_M$ learned by PPO

Initialize $s_0^i \sim \mu_0$, $i \in [N]$

for $t = 1$ to T **do**

for $i = 1$ to N **do**

 Choose $m(i) = \arg \min_{m \in [M]} |\frac{i}{N} - \frac{m}{M}|$

 Sample action $a_t^i \sim \tilde{\pi}^{\alpha_{m(i)}}(\cdot | s_t^i)$, observe reward r_t^i and new state s_{t+1}^i

end for

end for

311 **4. Proofs of Main Results**

312 In this section, we will provide proofs of Theorems 3.7-3.9.

313 *4.1. Proof of Theorem 3.7*

314 To prove Theorem 3.7, we need the following two Lemmas. We start by defining the
315 step state distribution $\mu_t^N := (\mu_t^{N,\alpha})_{\alpha \in \mathcal{I}}$ for notational simplicity

$$\mu_t^{N,\alpha}(\cdot) = \sum_{i \in \mathcal{V}_N} \delta_{s_t^i}(\cdot) \mathbf{1}_{\alpha \in (\frac{i-1}{N}, \frac{i}{N}]}. \quad (4.1)$$

316 Lemma 4.1 shows the convergence of the neighborhood empirical measure to the neigh-
317 borhood mean-field measure.

318 **Lemma 4.1** *Assume Assumptions 2.1, 3.3, 3.4 and 3.6. Under either condition (C1) or*
319 *(C2), for any policy ensemble $\pi \in \Pi$, we have*

$$\sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E}[\|\mu_t^{i,W_N} - \mu_t^{\alpha,W}\|_1] d\alpha \rightarrow 0, \quad \text{as } N \rightarrow \infty, \quad (4.2)$$

320 where $\mu_t^i = \mu_t^\alpha = \mu \in \mathcal{P}(\mathcal{S})$.

Moreover, if the graphon convergence in Assumption 2.1 is at rate $\mathcal{O}(\frac{1}{\sqrt{N}})$, then

$$\sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E}[\|\mu_t^{i,W_N} - \mu_t^{\alpha,W}\|_1] d\alpha = \mathcal{O}(\frac{1}{\sqrt{N}}).$$

321 **Proof of Lemma 4.1** We first prove (4.2) under the condition (C1) and then show (4.2)
322 also holds under the condition (C2).

323 **Case 1:** $\xi_{i,j}^N = W_N(\frac{i}{N}, \frac{j}{N})$. Note that under the condition (C1), $\mu_t^{i,W_N} = \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N,\beta} d\beta$
 324 by the definition of $\mu_t^{N,\alpha}$ in (4.1). Then

$$\begin{aligned}
 & \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} [\|\mu_t^{i,W_N} - \mu_t^{\alpha,W}\|_1] d\alpha \\
 &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \left[\left\| \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N,\beta} d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 \right] d\alpha \\
 &\leq \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \left[\left\| \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N,\beta} d\beta - \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{\beta} d\beta \right\|_1 \right] d\alpha \\
 &\quad + \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \left[\left\| \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{\beta} d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 \right] d\alpha \\
 &=: I_1 + I_2.
 \end{aligned}$$

325 For the term I_1 , we adapt Theorem 2 that works with local policy in [15] to our setting of
 326 global policy and have that under the policy ensemble π and N -agent policy $(\pi^1, \dots, \pi^N) :=$
 327 $\Gamma_N(\pi)$, with Γ_N defined in (3.7)

$$I_1 = \mathbb{E} \left[\left\| \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N,\beta} d\beta - \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{\beta} d\beta \right\|_1 \right] \rightarrow 0, \text{ as } N \rightarrow \infty.$$

328 Moreover, if the graphon convergence in Assumption 2.1 is at rate $\mathcal{O}(\frac{1}{\sqrt{N}})$, then the term
 329 I_1 is also at rate $\mathcal{O}(\frac{1}{\sqrt{N}})$.

330 By noting that $W_N(\alpha, \beta) = W_N(\frac{[N\alpha]}{N}, \frac{[N\beta]}{N})$,

$$\begin{aligned}
 I_2 &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left\| \int_{\mathcal{I}} W_N(\frac{[N\alpha]}{N}, \beta) \mu_t^{\beta} d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 d\alpha \\
 &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left\| \int_{\mathcal{I}} W_N(\alpha, \beta) \mu_t^{\beta} d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 d\alpha \\
 &= \int_{\mathcal{I}} \left\| \int_{\mathcal{I}} W_N(\alpha, \beta) \mu_t^{\beta} d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 d\alpha \\
 &= \sum_{s \in \mathcal{S}} \int_{\mathcal{I}} \left| \int_{\mathcal{I}} W_N(\alpha, \beta) \mu_t^{\beta}(s) d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta}(s) d\beta \right| d\alpha \\
 &\rightarrow 0,
 \end{aligned}$$

where the last inequality is from the fact in [39] that the convergence of $\|W_N - W\|_{\square} \rightarrow 0$ is equivalent to the convergence of

$$\|W_N - W\|_{L_{\infty} \rightarrow L_1} := \sup_{\|g\|_{\infty} \leq 1} \int_{\mathcal{I}} \left| \int_{\mathcal{I}} (W_N(\alpha, \beta) - W(\alpha, \beta)) g(\beta) d\beta \right| d\alpha \rightarrow 0.$$

331 Combining I_1 and I_2 , we prove (4.2) under the condition (C1).

332 **Case 2:** $\xi_{i,j}^N$ are random variables with Bernoulli($W_N(\frac{i}{N}, \frac{j}{N})$).

$$\begin{aligned}
 & \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \|\mu_t^{i, W_N} - \mu_t^{\alpha, W}\|_1 d\alpha \\
 = & \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \xi_{ij}^N \delta_{s_t^j} - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 d\alpha \\
 \leq & \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \xi_{ij}^N \delta_{s_t^j} - \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N, \beta} d\beta \right\|_1 d\alpha \\
 & + \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \left\| \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N, \beta} d\beta - \int_{\mathcal{I}} W(\alpha, \beta) \mu_t^{\beta} d\beta \right\|_1 d\alpha \\
 =: & I_1 + I_2.
 \end{aligned}$$

333 Note from **Case 1** that $I_2 \rightarrow 0$ as $N \rightarrow \infty$ and $I_2 = \mathcal{O}(\frac{1}{\sqrt{N}})$ if the graphon convergence in
 334 Assumption 2.1 is at rate $\mathcal{O}(\frac{1}{\sqrt{N}})$. Therefore, it is enough to estimate I_1 .

$$\begin{aligned}
 I_1 &= \mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \xi_{ij}^N \delta_{s_t^j} - \int_{\mathcal{I}} W_N(\frac{i}{N}, \beta) \mu_t^{N, \beta} d\beta \right\|_1 \\
 &\leq \mathbb{E} \left[\mathbb{E} \left[\sup_{f: \mathcal{S} \rightarrow \{-1, 1\}} \left\{ \frac{1}{N} \sum_{j=1}^N \xi_{ij}^N f(s_t^j) - \frac{1}{N} \sum_{j=1}^N W_N(\frac{i}{N}, \frac{j}{N}) f(s_t^j) \right\} \middle| s_t^1, \dots, s_t^N \right] \right].
 \end{aligned}$$

335 We proceed the same argument as in the proof of Theorem 6.3 in [23]. Precisely, conditioned
 336 on s_t^1, \dots, s_t^N , $\left\{ \xi_{ij}^N f(s_t^j) - W_N(\frac{i}{N}, \frac{j}{N}) f(s_t^j) \right\}_{j=1}^N$ is a sequence of independent mean-zero
 337 random variables bounded in $[-1, 1]$ due to $\mathbb{E}[\xi_{i,j}^N] = W_N(\frac{i}{N}, \frac{j}{N})$. This implies that each
 338 $\xi_{ij}^N f(s_t^j) - W_N(\frac{i}{N}, \frac{j}{N}) f(s_t^j)$ is a sub-Gaussian with variance bounded by 4. As a result,
 339 conditioned on s_t^1, \dots, s_t^N , $\left\{ \frac{1}{N} \sum_{j=1}^N \xi_{ij}^N f(s_t^j) - \frac{1}{N} \sum_{j=1}^N W_N(\frac{i}{N}, \frac{j}{N}) f(s_t^j) \right\}_{i=1}^N$ is a mean-zero
 340 sub-Gaussian random variable with variance $\frac{4}{N}$. By the equation (2.66) in [51], we have

$$\begin{aligned}
 I_1 &\leq \mathbb{E} \left[\mathbb{E} \left[\sup_{f: \mathcal{S} \rightarrow \{-1, 1\}} \left\{ \frac{1}{N} \sum_{j=1}^N \xi_{ij}^N f(s_t^j) - \frac{1}{N} \sum_{j=1}^N W_N(\frac{i}{N}, \frac{j}{N}) f(s_t^j) \right\} \middle| s_t^1, \dots, s_t^N \right] \right] \\
 &\leq \frac{\sqrt{8 \ln(2) |\mathcal{S}|}}{\sqrt{N}}.
 \end{aligned}$$

341 Therefore, combining I_1 and I_2 in **Case 2**, we show that when $\xi_{i,j}^N$ are random variables
 342 with Bernoulli($W_N(\frac{i}{N}, \frac{j}{N})$), (4.2) holds under the condition (C2). \square

343 Lemma 4.2 shows the convergence of the state distribution of N -agent game to the state
 344 distribution of GMFC.

345 **Lemma 4.2** Assume Assumptions 2.1, 3.3, 3.4 and 3.6. For any uniformly bounded family
 346 \mathcal{G} of functions $g^\alpha : \mathcal{S} \rightarrow \mathbb{R}$, we have

$$\sup_{\{g^\alpha\}_{\alpha \in \mathcal{I}} \in \mathcal{G}} \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[g^\alpha(s_t^i) - g^\alpha(s_t^\alpha)] \right| d\alpha \rightarrow 0, \quad (4.3)$$

where $s_0^i \sim \mu_0$, $s_0^\alpha \sim \mu_0$. Moreover, if the graphon convergence in Assumption 2.1 is at rate $\mathcal{O}(\frac{1}{\sqrt{N}})$, then

$$\sup_{\{g^\alpha\}_{\alpha \in \mathcal{I}} \in \mathcal{G}} \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[g^\alpha(s_t^i) - g^\alpha(s_t^\alpha)] \right| d\alpha = \mathcal{O}\left(\frac{1}{\sqrt{N}}\right).$$

347 **Proof of Lemma 4.2** The proof is by induction as follows. To do this, first introduce

$$l_{g^\alpha}^\beta(s, \mu, \pi) := \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} g^\alpha(s') P^\beta(s' | s, \mu, a) \pi(a | s, \mu).$$

348 (4.3) holds obviously at $t = 0$. Suppose that (4.3) holds at t . Then for any uniformly
 349 bounded function g^α with $|g^\alpha| \leq M_g$ at $t + 1$

$$\begin{aligned} & \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[g^\alpha(s_{t+1}^i) - g^\alpha(s_{t+1}^\alpha)] \right| d\alpha \\ &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^{\frac{i}{N}}(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^\alpha)] \right| d\alpha \\ &\leq \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^{\frac{i}{N}}(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] \right| d\alpha \\ &\quad + \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^i)] \right| d\alpha \\ &\quad + \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^\alpha)] \right| d\alpha \\ &= : I + II + III, \end{aligned} \quad (4.4)$$

350 where the first equality is by the law of total expectation.

First term of (4.4).

$$\begin{aligned} I &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^{\frac{i}{N}}(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] \right| d\alpha \\ &\leq M_g \left(L_P \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E}[\|\mu_t^{i, W_N} - \mu_t^{\alpha, W}\|_1] d\alpha + \tilde{L}_P \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \alpha - \frac{i}{N} \right| d\alpha \right) \\ &\rightarrow 0, \quad \text{as } N \rightarrow \infty \end{aligned}$$

351 where the second inequality is from the continuity of P , and the last inequality is from
 352 Lemma 4.1.

353 *Second term of (4.4).* One can view $l_{g^\alpha}^\alpha(s, \mu_t^{\alpha, W}, \pi^i)$ as a function of $s \in \mathcal{S}$ for any fixed
 354 $\mu_t^{\alpha, W}$ and π^i , $\alpha \in \mathcal{I}$. Note that $|l_{g^\alpha}^\alpha(s, \mu_t^{\alpha, W}, \pi^i)| \leq M_g$, where M_g is a constant independent
 355 of $\mu_t^{\alpha, W}$, π^i . Since (4.3) holds at t , then

$$\begin{aligned} II &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^i)] \right| d\alpha \\ &\rightarrow 0, \quad \text{as } N \rightarrow \infty. \end{aligned}$$

Third term of (4.4).

$$\begin{aligned} III &= \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[l_{g^\alpha}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^\alpha)] \right| d\alpha \\ &\leq M_g \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E}[\|\pi^i(s_t^\alpha) - \pi^\alpha(s_t^\alpha)\|_1] d\alpha \\ &\leq M_g L_\Pi \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \max_{\alpha \in (\frac{i-1}{N}, \frac{i}{N}]} \left| \frac{i}{N} - \alpha \right| d\alpha \\ &= \mathcal{O}\left(\frac{1}{N}\right), \end{aligned}$$

356 where the second inequality is by the uniform boundedness of g and the third inequality is
 357 from Assumption 3.6. \square

358 Now we are ready to prove Theorem 3.7. We start by defining \widehat{r}^α the aggregated reward
 359 over all possible actions under the policy π

$$\widehat{r}^\alpha(s, \mu, \pi) := \sum_{a \in \mathcal{A}} r^\alpha(s, \mu, a) \pi(a|s, \mu).$$

Proof of Theorem 3.7

$$\begin{aligned}
 & |V_N(\mu) - V(\mu)| \\
 &= \left| \sup_{\boldsymbol{\pi} \in \Pi^N} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^i(s_t^i, \mu_t^{i, W_N}, a_t^i) \right] - \sup_{\boldsymbol{\pi} \in \Pi} \int_{\mathcal{I}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, a_t^\alpha) \right] d\alpha \right| \\
 &\leq \sup_{\boldsymbol{\pi} \in \Pi} \left| \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^i(s_t^i, \mu_t^{i, W_N}, a_t^i) \right] - \int_{\mathcal{I}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, a_t^\alpha) \right] d\alpha \right| \\
 &= \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^i(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[\widehat{r}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^\alpha)] \right) d\alpha \right| \\
 &\leq \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^i(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[\widehat{r}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] \right) d\alpha \right| \\
 &\quad + \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[\widehat{r}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^i)] \right) d\alpha \right| \\
 &\quad + \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[\widehat{r}^\alpha(s_t^\alpha, \mu_t^{\alpha, W}, \pi^\alpha)] \right) d\alpha \right| \\
 &:= I + II + III, \tag{4.5}
 \end{aligned}$$

360 where we use (3.8) in the second inequality.

First term of (4.5).

$$\begin{aligned}
 I &= \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^i(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[\widehat{r}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] \right) d\alpha \right| \\
 &= \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^i(s_t^i, \mu_t^{i, W_N}, \pi^i)] - \mathbb{E}[\widehat{r}^i(s_t^i, \mu_t^{\alpha, W}, \pi^i)] \right) d\alpha \right| \\
 &\quad + \sup_{\boldsymbol{\pi} \in \Pi} \left| \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left(\mathbb{E}[\widehat{r}^i(s_t^i, \mu_t^{\alpha, W}, \pi^i)] - \mathbb{E}[\widehat{r}^\alpha(s_t^i, \mu_t^{\alpha, W}, \pi^i)] \right) d\alpha \right| \\
 &\leq \sup_{\boldsymbol{\pi}} L_r \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \mathbb{E} \|\mu_t^{i, W_N} - \mu_t^{\alpha, W}\|_1 d\alpha + \sup_{\boldsymbol{\pi}} \tilde{L}_r \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \frac{i}{N} - \alpha \right| d\alpha \\
 &= \mathcal{O}\left(\frac{1}{\sqrt{N}}\right), \tag{4.6}
 \end{aligned}$$

361 where the last equality is from Lemma 4.1 when the convergence in Assumption 2.1 is at
 362 rate $\mathcal{O}(1/\sqrt{N})$.

363 Second term of (4.5). From Lemma 4.2, we have $II = \mathcal{O}(\frac{1}{\sqrt{N}})$.

Third term of (4.5).

$$\begin{aligned}
 III &\leq \sup_{\pi} L_r \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \max_{s \in \mathcal{S}} \|\pi^i(s) - \pi^\alpha(s)\|_1 d\alpha \\
 &\leq L_r \tilde{L}_\Pi \sup_{\pi} \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \int_{(\frac{i-1}{N}, \frac{i}{N}]} \left| \frac{i}{N} - \alpha \right| d\alpha \\
 &= \mathcal{O}\left(\frac{1}{N}\right).
 \end{aligned}$$

364 Therefore, combining I, II and III yields the desired result. \square

365 4.2. Proof of Theorem 3.8

366 First, we see that (3.11) corresponds to the following optimal control problem

$$\begin{aligned}
 \tilde{V}^M(\tilde{\mu}) &:= \sup_{\tilde{\pi} \in \tilde{\Pi}_M} \tilde{J}^M(\tilde{\mu}, \tilde{\pi}) \\
 &= \sup_{\tilde{\pi} \in \tilde{\Pi}_M} \frac{1}{M} \sum_{m=1}^M \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(\tilde{s}_t^{\alpha_m}, \tilde{\mu}_t^{\alpha_m, W}, \tilde{a}_t^{\alpha_m}) \mid \tilde{s}_0^{\alpha_m} \sim \tilde{\mu}^{\alpha_m}, \tilde{a}_t^{\alpha_m} \sim \tilde{\pi}^{\alpha_m}(\cdot \mid \tilde{s}_t^{\alpha_m}) \right] \quad (4.7)
 \end{aligned}$$

367 The associated Q function of (4.7) is defined as

$$\begin{aligned}
 \tilde{Q}(\tilde{\mu}, \tilde{\pi}) &= \sup_{\tilde{\pi}'} \frac{1}{M} \sum_{m=1}^M \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(\tilde{s}_t^{\alpha_m}, \tilde{\mu}_t^{\alpha_m, W}, \tilde{a}_t^{\alpha_m}) \mid \tilde{s}_0^{\alpha_m} \sim \tilde{\mu}^{\alpha_m}, \tilde{a}_0^{\alpha_m} \sim \tilde{\pi}^{\alpha_m}(\cdot \mid \tilde{s}_t^{\alpha_m}) \right] \\
 &= R(\tilde{\mu}, \tilde{\pi}) + \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \sum_{t=1}^{\infty} \gamma^t \tilde{R}(\tilde{\mu}_t, \tilde{\pi}'), \quad (4.8)
 \end{aligned}$$

368 subject to $\tilde{\mu}_{t+1} = \tilde{\Phi}(\tilde{\mu}_t, \tilde{\pi})$, $\tilde{\mu}_0 = \tilde{\mu}$.

369 We first show the verification result and then prove the continuity property of \tilde{Q} in (4.8),
 370 which thus leads to Theorem 3.8.

371 **Lemma 4.3 (Verification)** *Assume Assumption 3.5. Then \tilde{Q} in (4.8) is the unique func-*
 372 *tion satisfying the Bellman equation (3.11). Furthermore, if there exists $\tilde{\pi}^* \in \arg \max_{\tilde{\Pi}_M} \tilde{Q}(\tilde{\mu}, \tilde{\pi})$*
 373 *for each $\tilde{\mu} \in \tilde{\mathcal{M}}_M$, then $\tilde{\pi}^* \in \tilde{\Pi}_M$ is an optimal stationary policy ensemble.*

374 The proof of Lemma 4.3 is standard and very similar to the proof of Proposition 3.3 in
 375 [23].

376 **Proof of Lemma 4.3** First, define $\frac{M_r}{1-\gamma}$ -bounded function space $\mathcal{Q} := \{f : \tilde{\mathcal{M}}_M \times \tilde{\Pi}_M \rightarrow$
 377 $[-\frac{M_r}{1-\gamma}, \frac{M_r}{1-\gamma}]\}$. Then we define a Bellman operator $B : \mathcal{Q} \rightarrow \mathcal{Q}$

$$(Bq)(\tilde{\mu}, \tilde{\pi}) := \tilde{R}(\tilde{\mu}, \tilde{\pi}) + \gamma \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} q(\tilde{\Phi}(\tilde{\mu}, \tilde{\pi}), \tilde{\pi}'),$$

378 One can show that B is a contraction operator with the module- γ . By Banach fixed point
 379 theorem, B admits a unique fixed point. As \tilde{Q} function of (4.8) satisfies $B\tilde{Q} = \tilde{Q}$, \tilde{Q} is
 380 unique solution of (3.11).

381 We next define $B^{\tilde{\pi}'}$: $\mathcal{Q} \rightarrow \mathcal{Q}$ under the policy ensemble $\tilde{\pi}' \in \tilde{\Pi}_M$ with

$$(B^{\tilde{\pi}'})_q(\tilde{\mu}, \tilde{\pi}) := \tilde{R}(\tilde{\mu}, \tilde{\pi}) + \gamma q(\tilde{\Phi}(\tilde{\mu}, \tilde{\pi}), \tilde{\pi}').$$

382 Similarly, we can show that $B^{\tilde{\pi}'}$ is a contraction map with the module- γ and thus admits a
 383 unique fixed point, which is denoted as $\tilde{Q}^{\tilde{\pi}'}$. From this, we have

$$\begin{aligned} \tilde{Q}^{\tilde{\pi}^*}(\tilde{\mu}, \tilde{\pi}) &= \tilde{R}(\tilde{\mu}, \tilde{\pi}) + \gamma \tilde{Q}^{\tilde{\pi}^*}(\tilde{\Phi}(\tilde{\mu}, \tilde{\pi}), \tilde{\pi}^*) \\ &= \tilde{R}(\tilde{\mu}, \tilde{\pi}) + \gamma \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \tilde{Q}(\tilde{\Phi}(\tilde{\mu}, \tilde{\pi}), \tilde{\pi}') = \tilde{Q}(\tilde{\mu}, \tilde{\pi}), \end{aligned}$$

384 which implies $\tilde{\pi}^*$ is an optimal policy ensemble. □

385 **Lemma 4.4** *Let Assumptions 3.4, 3.5 hold. Assume further $\gamma \cdot (1 + L_P + L_\Pi) < 1$. Then*
 386 *\tilde{Q} in (4.8) is continuous.*

387 **Proof of Lemma 4.4** We will show that as $\tilde{\mu}_n \rightarrow \tilde{\mu}$, $\tilde{\pi}_n \rightarrow \tilde{\pi}$ in the sense that $\frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}^{\alpha_m} -$
 388 $\tilde{\mu}_n^{\alpha_m}\|_1 + \frac{1}{M} \sum_{m=0}^{M-1} \max_{s \in \mathcal{S}} \|\tilde{\pi}^{\alpha_m}(\tilde{\mu}^{\alpha_m, W}) - \tilde{\pi}_n^{\alpha_m}(\tilde{\mu}_n^{\alpha_m, W})\|_1 \rightarrow 0$,

$$\tilde{Q}(\tilde{\mu}_n, \tilde{\pi}_n(\tilde{\mu}_n)) \rightarrow \tilde{Q}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu})).$$

389 From (4.8) and (3.13),

$$\begin{aligned} &|\tilde{Q}(\tilde{\mu}_n, \tilde{\pi}_n) - \tilde{Q}(\tilde{\mu}, \tilde{\pi})| \\ &\leq \left| \tilde{R}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu})) + \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \sum_{t=1}^{\infty} \gamma^t \tilde{R}(\tilde{\mu}_t, \tilde{\pi}'(\tilde{\mu}_t)) - \tilde{R}(\tilde{\mu}_n, \tilde{\pi}_n(\tilde{\mu}_n)) + \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \sum_{t=1}^{\infty} \gamma^t \tilde{R}(\tilde{\mu}_{n,t}, \tilde{\pi}'(\tilde{\mu}_{n,t})) \right| \\ &\leq \left| \tilde{R}(\tilde{\mu}, \tilde{\pi}(\tilde{\mu})) - \tilde{R}(\tilde{\mu}_n, \tilde{\pi}_n(\tilde{\mu}_n)) \right| + \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \sum_{t=1}^{\infty} \gamma^t \left| \tilde{R}(\tilde{\mu}_{n,t}, \tilde{\pi}'(\tilde{\mu}_{n,t})) - \tilde{R}(\tilde{\mu}_t, \tilde{\pi}'(\tilde{\mu}_t)) \right| \\ &\leq L_r \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}^{\alpha_m, W} - \tilde{\mu}_n^{\alpha_m, W}\|_1 + M_r \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}^{\alpha_m} - \tilde{\mu}_n^{\alpha_m}\|_1 d\alpha \\ &\quad + M_r \cdot \frac{1}{M} \sum_{m=0}^{M-1} \max_{s \in \mathcal{S}} \|\tilde{\pi}^{\alpha}(\tilde{\mu}^{\alpha_m, W}) - \tilde{\pi}_n^{\alpha}(\tilde{\mu}_n^{\alpha_m, W})\|_1 d\alpha \\ &\quad + \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \sum_{t=1}^{\infty} \gamma^t \cdot \left((L_r + L_\Pi) \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}_t^{\alpha_m, W} - \tilde{\mu}_{n,t}^{\alpha_m, W}\|_1 + M_r \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}_t^{\alpha_m} - \tilde{\mu}_{n,t}^{\alpha_m}\|_1 \right) \\ &\leq (L_r + M_r) \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}^{\alpha_m} - \tilde{\mu}_n^{\alpha_m}\|_1 + M_r \cdot \frac{1}{M} \sum_{m=0}^{M-1} \max_{s \in \mathcal{S}} \|\tilde{\pi}^{\alpha}(\tilde{\mu}^{\alpha_m, W}) - \tilde{\pi}_n^{\alpha}(\tilde{\mu}_n^{\alpha_m, W})\|_1 d\alpha \\ &\quad + \sup_{\tilde{\pi}' \in \tilde{\Pi}_M} \sum_{t=1}^{\infty} \gamma^t \cdot (L_r + L_\Pi + M_r) \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}_t^{\alpha_m} - \tilde{\mu}_{n,t}^{\alpha_m}\|_1. \end{aligned}$$

390 By induction, we obtain

$$\begin{aligned}
 & \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}_t^{\alpha_m} - \tilde{\mu}_{n,t}^{\alpha_m}\|_1 \\
 = & \frac{1}{M} \sum_{m=0}^{M-1} \sum_{s' \in \mathcal{S}} \left| \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} P^{\alpha_m}(s'|s, a, \tilde{\mu}_{t-1}^{\alpha_m, W}) \tilde{\mu}_{t-1}^{\alpha_m}(s) \tilde{\pi}(\alpha|s, \tilde{\mu}_{t-1}^{\alpha_m, W}) \right. \\
 & \left. - \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} P^{\alpha_m}(s'|s, a, \tilde{\mu}_{n,t-1}^{\alpha_m, W}) \tilde{\mu}_{n,t-1}^{\alpha_m}(s) \tilde{\pi}(\alpha|s, \tilde{\mu}_{n,t-1}^{\alpha_m, W}) \right| \\
 \leq & (L_P + L_\Pi + 1) \cdot \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}_{t-1}^{\alpha} - \tilde{\mu}_{n,t-1}^{\alpha}\|_1 \\
 \leq & \dots \leq (L_P + L_\Pi + 1)^{(t-1)} \frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}_1^{\alpha} - \tilde{\mu}_{n,1}^{\alpha}\|_1.
 \end{aligned}$$

391 Therefore, if $\gamma \cdot (1 + L_P + L_\Pi) < 1$, then

$$|\tilde{Q}(\tilde{\mu}_n, \tilde{\pi}_n) - \tilde{Q}(\tilde{\mu}, \tilde{\pi})| \leq C \left(\frac{1}{M} \sum_{m=0}^{M-1} \|\tilde{\mu}^{\alpha_m} - \tilde{\mu}_n^{\alpha_m}\|_1 + \frac{1}{M} \sum_{m=0}^{M-1} \max_{s \in \mathcal{S}} \|\tilde{\pi}^{\alpha_m}(\tilde{\mu}^{\alpha_m, W}) - \tilde{\pi}_n^{\alpha_m}(\tilde{\mu}_n^{\alpha_m, W})\|_1 \right).$$

392 where C is a constant depending on L_r, M_r, L_P, L_Π . \square

393 Now we prove Theorem 3.8.

394 **Proof of Theorem 3.8** By Lemma 4.4, along with the compactness of $\tilde{\Pi}_M$, there exists
 395 $\tilde{\pi}^* \in \tilde{\Pi}_M$ such that $\tilde{\pi}^* \in \arg \max_{\tilde{\pi} \in \tilde{\Pi}_M} Q(\tilde{\mu}, \tilde{\pi})$. By Lemma 4.3, there exists an optimal policy

396 ensemble $\tilde{\pi}^* \in \tilde{\Pi}_M$. \square

397 4.3. Proof of Theorem 3.9

398 We first prove the following Lemma, which shows that GMFC and block GMFC become
 399 increasingly close to each other as the number of blocks becomes larger.

400 **Lemma 4.5** Under Assumptions 3.3, 3.4 and 3.6, we have

$$\begin{aligned}
 \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}] } \|\mu_t^{\alpha, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 d\alpha & \leq \left[(1 + L_P + L_\Pi)^t - 1 \right] \frac{\tilde{L}_\Pi + \tilde{L}_P + 2(L_P + L_\Pi)L_W}{M} + \frac{2L_W}{M}, \\
 \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}] } \|\mu_t^{\alpha} - \tilde{\mu}_t^{\alpha_m}\|_1 d\alpha & \leq \left[(1 + L_P + L_\Pi)^t - 1 \right] \frac{\tilde{L}_\Pi + \tilde{L}_P + 2(L_P + L_\Pi)L_W}{M}.
 \end{aligned}$$

Proof of Lemma 4.5

$$\begin{aligned}
 & \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^{\alpha, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 d\alpha \quad (4.9) \\
 \leq & \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^{\alpha, W} - \mu_t^{\alpha_m, W}\|_1 d\alpha + \frac{1}{M} \sum_{m=1}^M \|\mu_t^{\alpha_m, W} - \bar{\mu}_t^{\alpha_m, W}\|_1 \\
 & + \frac{1}{M} \sum_{m=1}^M \|\bar{\mu}_t^{\alpha_m, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1,
 \end{aligned}$$

401 where $\bar{\mu}^{\alpha_m, W} := \frac{1}{M} \sum_{m'=1}^M W(\alpha_m, \alpha_{m'}) \mu^{\alpha_{m'}}$.

402 By the definition of $\mu_t^{\alpha, W}$, $\mu_t^{\alpha_m, W}$ in (2.10), $\tilde{\mu}_t^{\alpha_m, W}$ in (3.12) and $\bar{\mu}^{\alpha_m, W}$, together with the
 403 Lipschitz continuity of W in Assumption 3.3,

$$\sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^{\alpha, W} - \mu_t^{\alpha_m, W}\|_1 d\alpha \leq \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^\alpha - \mu_t^{\alpha_m}\|_1 d\alpha + \frac{L_W}{M}, \quad (4.10)$$

$$\frac{1}{M} \sum_{m=1}^M \|\mu_t^{\alpha_m, W} - \bar{\mu}_t^{\alpha_m, W}\|_1 \leq \frac{L_W}{M}, \quad (4.11)$$

$$\frac{1}{M} \sum_{m=1}^M \|\bar{\mu}_t^{\alpha_m, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 \leq \frac{1}{M} \sum_{m=1}^M \|\mu_t^{\alpha_m} - \tilde{\mu}_t^{\alpha_m}\|_1. \quad (4.12)$$

404 Plugging these into (4.9),

$$\sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^{\alpha, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 d\alpha \leq A_t + \frac{2L_W}{M}, \quad (4.13)$$

405 where $A_t := \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^\alpha - \mu_t^{\alpha_m}\|_1 d\alpha + \frac{1}{M} \sum_{m=1}^M \|\mu_t^{\alpha_m} - \tilde{\mu}_t^{\alpha_m}\|_1$.

406 On the other hand,

$$\sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^\alpha - \tilde{\mu}_t^{\alpha_m}\|_1 d\alpha \quad (4.14)$$

$$\leq \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^\alpha - \mu_t^{\alpha_m}\|_1 d\alpha + \frac{1}{M} \sum_{m=1}^M \|\mu_t^{\alpha_m} - \tilde{\mu}_t^{\alpha_m}\|_1 = A_t. \quad (4.15)$$

407 Therefore, it is enough to estimate A_t . We next estimate A_{t+1} by an inductive way. Note

408 that $A_0 = 0$.

$$\begin{aligned}
 & A_{t+1} \\
 = & \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_{t+1}^\alpha - \mu_{t+1}^{\alpha_m}\|_1 d\alpha + \frac{1}{M} \sum_{m=1}^M \|\mu_{t+1}^{\alpha_m} - \tilde{\mu}_{t+1}^{\alpha_m}\|_1 \\
 = & \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \left\| \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left(P^\alpha(\cdot|s, \mu_t^{\alpha, W}, a) \pi^\alpha(a|s, \mu_t^{\alpha, W}) \mu_t^\alpha(s) \right. \right. \\
 & \left. \left. - P^{\alpha_m}(\cdot|s, a, \mu_t^{\alpha_m, W}) \mu_t^{\alpha_m}(s) \pi^{\alpha_m}(a|s, \mu_t^{\alpha_m, W}) \right) \right\|_1 d\alpha \\
 & + \frac{1}{M} \sum_{m=1}^M \left\| \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left(P^{\alpha_m}(\cdot|s, \mu_t^{\alpha_m, W}, a) \pi^{\alpha_m}(a|s, \mu_t^{\alpha_m, W}) \mu_t^{\alpha_m}(s) \right. \right. \\
 & \left. \left. - P^{\alpha_m}(\cdot|s, a, \tilde{\mu}_t^{\alpha_m, W}) \tilde{\mu}_t^{\alpha_m}(s) \tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}_t^{\alpha_m, W}) \right) \right\|_1 \\
 \leq & \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \left((L_P + L_\Pi) \cdot \|\mu_t^{\alpha, W} - \mu_t^{\alpha_m, W}\|_1 + \frac{\tilde{L}_\Pi}{M} + \|\mu_t^\alpha - \mu_t^{\alpha_m}\|_1 \right) d\alpha \\
 & + \frac{1}{M} \sum_{m=1}^M \left((L_P + L_\Pi) \cdot \|\mu_t^{\alpha_m, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 + \frac{\tilde{L}_P}{M} + \|\mu_t^{\alpha_m} - \tilde{\mu}_t^{\alpha_m}\|_1 \right) \\
 \leq & (1 + L_P + L_\Pi) A_t + (\tilde{L}_\Pi + \tilde{L}_P + 2(L_P + L_\Pi)L_W) \frac{1}{M},
 \end{aligned}$$

409 where the second equality is from (3.4) and (3.14), and we use Assumptions 3.3, 3.4 and 3.6
 410 in the third inequality, and we use (4.10)-(4.12) in the last inequality.
 411 By induction, we have

$$A_{t+1} \leq \left[(1 + L_P + L_\Pi)^t - 1 \right] \frac{\tilde{L}_\Pi + \tilde{L}_P + 2(L_P + L_\Pi)L_W}{M}.$$

412

□

413 Based on Lemma 4.5, we have the following Proposition.

414 **Proposition 4.6** Assume Assumptions 3.3, 3.4, 3.5, 3.6, and $\gamma \cdot (L_P + L_\Pi + 1) < 1$. Then
 415 we have for any $\mu \in \mathcal{P}(\mathcal{S})$

$$\sup_{\boldsymbol{\pi} \in \Pi} |\tilde{J}^M(\mu, \boldsymbol{\pi}) - J(\mu, \boldsymbol{\pi})| \rightarrow 0, \quad \text{as } M \rightarrow +\infty, \quad (4.16)$$

416 where \tilde{J}^M and J are given in (4.7) and (2.11), respectively.

417 **Proof of Proposition 4.6** Recall from (3.12) that

$$\tilde{J}^M(\mu, \tilde{\boldsymbol{\pi}}) = \sum_{t=0}^{\infty} \gamma^t \tilde{R}(\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\pi}}(\tilde{\boldsymbol{\mu}}_t)),$$

418 subject to $\tilde{\mu}_{t+1}^{\alpha_m} = \tilde{\Phi}^{\alpha_m}(\tilde{\mu}_t^{\alpha_m}, \tilde{\pi}^{\alpha_m})$, $t \in \mathbb{N}_+$, $\tilde{\mu}_0^\alpha \equiv \mu$, and $\tilde{\mu}_t^{\alpha_m, W}$ given in (3.12).

$$J(\mu, \boldsymbol{\pi}) = \sum_{t=0}^{\infty} \gamma^t R(\boldsymbol{\mu}_t, \boldsymbol{\pi}(\boldsymbol{\mu}_t)),$$

419 subject to $\mu_{t+1}^\alpha = \Phi^\alpha(\mu_t^\alpha, \pi^\alpha)$, $t \in \mathbb{N}_+$, $\mu_0^\alpha \equiv \mu$, and $\mu_t^{\alpha, W}$ given in (2.10). Since $\tilde{\boldsymbol{\pi}} :=$
 420 $(\tilde{\pi}^{\alpha_m})_{m \in [M]} \in \tilde{\boldsymbol{\Pi}}_M$ can be viewed as a piecewise-constant projection of $\boldsymbol{\pi} \in \boldsymbol{\Pi}$ onto $\tilde{\boldsymbol{\Pi}}_M$.
 421 Then,

$$\begin{aligned} & \sup_{\boldsymbol{\pi} \in \boldsymbol{\Pi}} |\tilde{J}^M(\mu, \boldsymbol{\pi}) - J(\mu, \boldsymbol{\pi})| \\ & \leq \sup_{\boldsymbol{\pi} \in \tilde{\boldsymbol{\Pi}}_M} \sum_{t=0}^{\infty} \gamma^t \left| \tilde{R}(\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\pi}}(\tilde{\boldsymbol{\mu}}_t)) - R(\boldsymbol{\mu}_t, \boldsymbol{\pi}(\boldsymbol{\mu}_t)) \right| \\ & \leq \sup_{\boldsymbol{\pi} \in \tilde{\boldsymbol{\Pi}}_M} \sum_{t=0}^{\infty} \gamma^t \left| \tilde{R}(\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\pi}}(\tilde{\boldsymbol{\mu}}_t)) - R(\boldsymbol{\mu}_t, \tilde{\boldsymbol{\pi}}(\boldsymbol{\mu}_t)) \right| + \sup_{\boldsymbol{\pi} \in \tilde{\boldsymbol{\Pi}}_M} \sum_{t=0}^{\infty} \gamma^t \left| R(\boldsymbol{\mu}_t, \tilde{\boldsymbol{\pi}}(\boldsymbol{\mu}_t)) - R(\boldsymbol{\mu}_t, \boldsymbol{\pi}(\boldsymbol{\mu}_t)) \right| \\ & := I + II. \end{aligned}$$

422 In terms of the term I , we first estimate $\left| \tilde{R}(\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\pi}}) - R(\boldsymbol{\mu}_t, \tilde{\boldsymbol{\pi}}) \right|$:

$$\begin{aligned} & \left| \tilde{R}(\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\pi}}(\tilde{\boldsymbol{\mu}}_t)) - R(\boldsymbol{\mu}_t, \tilde{\boldsymbol{\pi}}(\boldsymbol{\mu}_t)) \right| \\ & = \left| \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r^{\alpha_m}(s, a, \tilde{\mu}_t^{\alpha_m, W}) \tilde{\mu}_t^{\alpha_m}(s) \tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}_t^{\alpha_m, W}) d\alpha \right. \\ & \quad \left. - \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r^\alpha(s, a, \mu_t^{\alpha, W}) \mu_t^\alpha(s) \tilde{\pi}^{\alpha_m}(a|s, \mu_t^{\alpha, W}) d\alpha \right| \\ & \leq \left| \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} (r^{\alpha_m}(s, a, \tilde{\mu}_t^{\alpha_m, W}) - r^\alpha(s, a, \mu_t^{\alpha, W})) \tilde{\mu}_t^{\alpha_m}(s) \tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}_t^{\alpha_m, W}) d\alpha \right. \\ & \quad \left. + \left| \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r^\alpha(s, a, \mu_t^{\alpha, W}) (\tilde{\mu}_t^{\alpha_m}(s) - \mu_t^\alpha(s)) \tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}_t^{\alpha_m, W}) d\alpha \right| \right. \\ & \quad \left. + \left| \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r^\alpha(s, a, \mu_t^{\alpha, W}) \tilde{\mu}_t^\alpha(s) (\tilde{\pi}^{\alpha_m}(a|s, \tilde{\mu}_t^{\alpha_m, W}) - \tilde{\pi}^{\alpha_m}(a|s, \mu_t^{\alpha, W})) d\alpha \right| \right| \\ & \leq L_r \cdot \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \|\mu_t^{\alpha, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 d\alpha + \frac{\tilde{L}_r}{M} \\ & \quad + M_r \cdot \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \|\mu_t^\alpha - \tilde{\mu}_t^{\alpha_m}\|_1 d\alpha + M_r L_\Pi \cdot \sum_{m=1}^M \int_{\left(\frac{m-1}{M}, \frac{m}{M}\right]} \|\mu_t^{\alpha, W} - \tilde{\mu}_t^{\alpha_m, W}\|_1 d\alpha. \end{aligned}$$

423 By Lemma 4.5,

$$I \leq \frac{C(\gamma, L_\Pi, L_P, L_W, L_r, M_r)}{M}.$$

424 For the term II ,

$$\begin{aligned}
 & \sup_{\boldsymbol{\pi} \in \Pi} \sum_{t=0}^{\infty} \gamma^t \left| R(\boldsymbol{\mu}_t, \tilde{\boldsymbol{\pi}}(\boldsymbol{\mu}_t)) - R(\boldsymbol{\mu}_t, \boldsymbol{\pi}(\boldsymbol{\mu}_t)) \right| \\
 & \leq \sup_{\boldsymbol{\pi} \in \Pi} \sum_{t=0}^{\infty} \gamma^t M_r \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \max_{s \in \mathcal{S}} \|\pi^\alpha - \pi^{\alpha^m}\|_1 d\alpha \\
 & \quad + \sup_{\boldsymbol{\pi} \in \Pi} \sum_{t=0}^{\infty} \gamma^t M_r \sum_{m=1}^M \int_{(\frac{m-1}{M}, \frac{m}{M}]} \|\mu_t^{\alpha, W} - \tilde{\mu}_t^{\alpha^m, W}\|_1 d\alpha \\
 & \leq \frac{C(\gamma, L_\Pi, L_P, L_W, L_r, M_r)}{M}.
 \end{aligned}$$

425

□

426 **Proof of Theorem 3.9** Suppose that $\tilde{\boldsymbol{\pi}}^* \in \tilde{\Pi}_M \subset \Pi$ and $(\pi^{1,*}, \dots, \pi^{N,*}) \in \Pi^N$ are optimal
 427 policies of the problems (4.7) and (2.7), respectively. From Proposition 4.6, for any $\varepsilon > 0$,
 428 there exists sufficiently large $M_\varepsilon > 0$

$$|\tilde{J}^{M_\varepsilon}(\mu, \tilde{\boldsymbol{\pi}}^*) - J(\mu, \tilde{\boldsymbol{\pi}}^*)| \leq \frac{\varepsilon}{3},$$

429 where by (3.8), $\boldsymbol{\pi}^{N,*} := \sum_{i=1}^N \pi^{i,*} \mathbf{1}_{\alpha \in (\frac{i-1}{N}, \frac{i}{N}]}$

430 From Theorem 3.7, for any $\varepsilon > 0$, there exists N_ε such that for all $N \geq N_\varepsilon$

$$|J_N(\mu, \tilde{\pi}^{1,*}, \dots, \tilde{\pi}^{N,*}) - J(\mu, \tilde{\boldsymbol{\pi}}^*)| \leq \frac{\varepsilon}{3}, \quad |J_N(\mu, \pi^{1,*}, \dots, \pi^{N,*}) - J(\mu, \boldsymbol{\pi}^{N,*})| \leq \frac{\varepsilon}{3}.$$

431 Then we have

$$\begin{aligned}
 & J_N(\mu, \tilde{\pi}^{1,*}, \dots, \tilde{\pi}^{N,*}) - J_N(\mu, \pi^{1,*}, \dots, \pi^{N,*}) \\
 & \geq \underbrace{J_N(\mu, \tilde{\pi}^{1,*}, \dots, \tilde{\pi}^{N,*}) - J(\mu, \tilde{\boldsymbol{\pi}}^*)}_{I_1} + \underbrace{J(\mu, \tilde{\boldsymbol{\pi}}^*) - \tilde{J}^{M_\varepsilon}(\mu, \tilde{\boldsymbol{\pi}}^*)}_{I_2} \\
 & \quad + \underbrace{\tilde{J}^{M_\varepsilon}(\mu, \tilde{\boldsymbol{\pi}}^*) - \tilde{J}^{M_\varepsilon}(\mu, \boldsymbol{\pi}^{N,*})}_{I_3} + \underbrace{\tilde{J}^{M_\varepsilon}(\mu, \boldsymbol{\pi}^{N,*}) - J_N(\mu, \pi^{1,*}, \dots, \pi^{N,*})}_{I_4} \\
 & \geq -\frac{\varepsilon}{3} - \frac{\varepsilon}{3} - \frac{\varepsilon}{3} = -\varepsilon.
 \end{aligned}$$

432 where $I_3 \geq 0$ due to the optimality of $\tilde{\boldsymbol{\pi}}^*$ for $\tilde{V}^{M_\varepsilon}$. This means that the optimal policy of
 433 block GMFC provides an ε -optimal policy for the multi-agent system with $(\tilde{\pi}_1^*, \dots, \tilde{\pi}_N^*) :=$
 434 $\Gamma_N(\tilde{\boldsymbol{\pi}}^*)$. □

435 5. Experiments

436 In this section, we provide an empirical verification of our theoretical results, with two
 437 examples adapted from existing works on learning MFGs [16, 10] and learning GMFGs [15].

438 5.1. SIS Graphon Model

439 We consider a SIS graphon model in [16] under a cooperative setting. In this model,
 440 each agent $\alpha \in \mathcal{I}$ shares a state space $\mathcal{S} = \{S, I\}$ and an action space $\mathcal{A} = \{C, NC\}$, where
 441 S is susceptible, I is infected, C represents keeping contact with others, and NC means
 442 keeping social distance. The transition probability of each agent α is represented as follows

$$\begin{aligned} P^\alpha(s_{t+1} = I | s_t = S, a_t = C, \mu_t^{\alpha, W}) &= \beta_1 \mu_t^{\alpha, W}(I), \\ P^\alpha(s_{t+1} = I | s_t = S, a_t = NC, \mu_t^{\alpha, W}) &= \beta_2 \mu_t^{\alpha, W}(I), \\ P^\alpha(s_{t+1} = S | s_t = I, \mu_t^{\alpha, W}) &= \delta, \end{aligned}$$

443 where β_1 is the infection rate with keeping contact with others, β_2 is the infection rate
 444 under social distance, and δ is the fixed recovery rate. We assume $0 < \beta_2 < \beta_1$, meaning
 445 that keeping social distance can reduce the risk of being infected. The individual reward
 446 function is defined as

$$r^\alpha(s, \mu_t^{\alpha, W}, a) = -c_1 \mathbf{1}_{\{I\}}(s) - c_2 \mathbf{1}_{\{NC\}}(a) - c_3 \mathbf{1}_{\{I\}}(s) \mathbf{1}_{\{C\}}(a),$$

447 where c_1 represents the cost of being infected such as the cost of medical treatment, c_2
 448 represents the cost of keeping social distance, and c_3 represents the penalty of going out if
 449 the agent is infected.

450 In our experiment, we set $\beta_1=0.8$, $\beta_2=0$, $\delta = 0.3$ for the transition dynamics and $c_1=2$,
 451 $c_2=0.3$, $c_3 = 0.5$ for the reward function. The initial mean field μ_0 is taken as the uniform
 452 distribution. We set the episode length to 50.

453 5.2. Malware Spread Graphon Model

We consider a malware spread model in [10] under a cooperative setting. In this model,
 let $\mathcal{S} = \{0, 1, \dots, K-1\}$, $K \in \mathbb{N}$, denote the health level of the agent, where $s_t = 0$ and
 $s_t = K-1$ represents the best level and the worst level, respectively. All agents can take
 two actions: $a_t = 0$ means doing nothing, and $a_t = 1$ means repairing. The state transition
 is given by

$$s_{t+1} = \begin{cases} s_t + \lfloor (K - s_t) \chi_t \rfloor, & \text{if } a_t = 0, \\ 0, & \text{if } a_t = 1, \end{cases}$$

454 where $\chi_t, t \in \mathbb{N}$ are i.i.d. random variables with a certain probability distribution. Then
 455 after taking action a_t , agent α will receive an individual reward

$$r^\alpha(s_t, \mu_t^{\alpha, W}, a_t) = -(c_1 + \langle \mu_t^{\alpha, W} \rangle) s_t / K - c_2 a_t.$$

456 Here considering the heterogeneity of agents, we use $W(\alpha, \beta)$ to denote the *importance* effect
 457 of agent β on agent α . $\langle \mu_t^{\alpha, W} \rangle := \int_{\beta \in \mathcal{I}} \sum_{s \in \mathcal{S}} s W(\alpha, \beta) \mu_t^\beta(s) d\beta$ is the risk of being infected
 458 by other agents and c_2 is the cost of taking action a_t .

459 In our experiment, we set $K=3$, $c_1=0.3$, and $c_2=0.5$. In addition, to stabilize the training
 460 of the RL agent, we fix χ_t to a static value, i.e., 0.7. In this model, we set the episode length
 461 to 10.

462 5.3. Performance of N-agent GMFC on Multi-Agent System

463 For both models, we use PPO [47] to train the block GMFC agent in the infinite-agent
 464 environment and obtain the policy ensembles and further use Algorithm 2 to deploy them
 465 in the finite-agent environment. We test the performance of N-agent GMFC with 10 blocks
 466 to different numbers of agents, i.e., from 10 to 100. For each case, we run 1000 times of
 467 simulations and show the mean and standard variation (Green shadows in Figure 1 and
 468 Figure 2) of the mean episode reward. We can see that in both scenarios and for different
 469 types of graphons, the mean episode rewards of the N-agent GMFC become increasingly
 470 close to that of block GMFC as the number of agents grows. (See Figure 1 and Figure 2).
 471 This verifies our theoretical findings empirically.

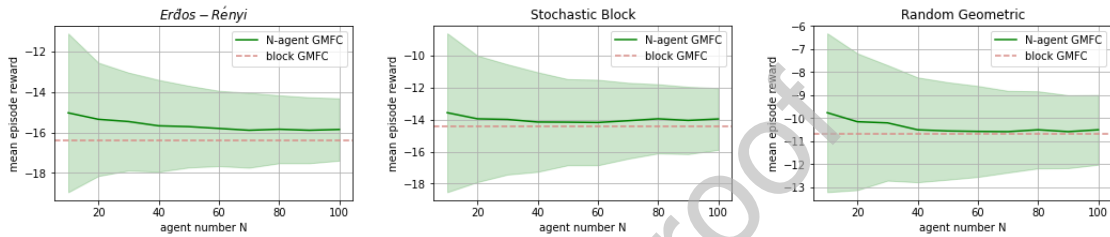


Figure 1: Experiments for different graphons in SIS finite-agent environment

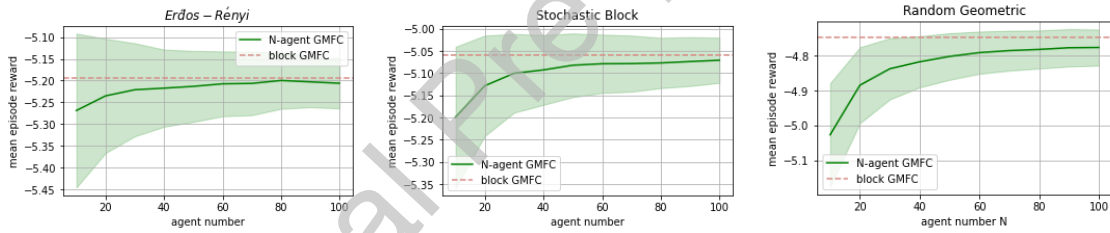


Figure 2: Experiments for different graphons in Malware Spread finite-agent environment

472 5.4. Comparison with Other Algorithms

473 For different types of graphons, we compare our algorithm N-agent GMFC with three
 474 existing MARL algorithms, including two independent learning algorithms, i.e., independent
 475 DQN [40], independent PPO [47] and a powerful centralized-training-and-decentralized-
 476 execution(CTDE)-based algorithm QMIX [46]. We test the performance of those algorithms
 477 with different numbers of blocks, i.e., 2, 5, 10, to the multi-agent systems with 40 agents.
 478 The results are reported in Table 1 and Table 2.

479 In the SIS graphon model, N-agent GMFC shows dominating performance in most cases
 480 and outperforms independent algorithms by a large margin. Only QMIX can reach compa-
 481 rable results. And in the malware spread graphon model, N-agent GMFC outperforms other
 482 algorithms in more than half of the cases. Only independent DQN has comparable perfor-
 483 mance in this environment. And we can see that in both environments, the performance

484 gap between N-agent GMFC and other MARL algorithms is shrinking as the number of
 485 blocks goes larger. This is mainly because the action space of block GMFC increases more
 486 quickly than MARL algorithms as the block number increases. And it is hard to train RL
 487 agents when the action space is too large.

488 Beyond the visible results shown in Tables 1 and 2, when the number of agents N grows
 489 larger, classic MARL methods become infeasible because of the curse of dimensionality
 490 and the restriction of memory storage, while N-agent GMFC is trained only once and
 491 independent of the number of agents N , hence is easier to scale up in a large-scale regime
 492 and enjoys a more stable performance. We can see that N-agent GMFC shows more stable
 493 results when N increases as shown in Figure 1 and Figure 2.

Table 1: Mean Episode Reward for SIS with 40 agents

Graphon Type	M	Algorithm			
		N-agent GMFC	I-DQN	I-PPO	QMIX
Erdős Rényi	2	-15.37	-17.58	-20.63	-20.51
	5	-15.74	-16.17	-20.42	16.94
	10	-15.67	-17.55	-21.38	-14.45
Stochastic Block	2	-13.58	-16.05	-18.38	-17.69
	5	-13.67	-15.91	-20.13	-13.79
	10	-13.57	-15.52	-14.87	-13.86
Random Geometric	2	-12.45	-17.93	-14.82	-14.52
	5	-9.82	-12.81	-12.99	-10.84
	10	-10.52	-11.68	-12.66	-12.60

Table 2: Mean Episode Reward for Malware Spread with 40 agents

Graphon Type	M	Algorithm			
		N-agent GMFC	I-DQN	I-PPO	QMIX
Erdős Rényi	2	-5.21	-5.11	-5.31	-6.05
	5	-5.21	-5.30	-5.26	-6.13
	10	-5.21	-5.14	-5.27	-5.21
Stochastic Block	2	-5.16	-5.21	-5.37	-5.88
	5	-5.10	-5.19	-5.31	-5.70
	10	-5.09	-5.05	-5.28	-5.27
Random Geometric	2	-5.02	-5.21	-5.27	-5.35
	5	-4.85	-5.03	-5.04	-5.05
	10	-4.82	-4.83	-5.14	-4.83

494 5.5. Implementation Details

495 We use three graphons in our experiments: (1) Erdős Rényi: $W(\alpha, \beta) = 0.8$; (2) Stochastic
 496 block model: $W(\alpha, \beta) = 0.9$, if $0 \leq \alpha, \beta \leq 0.5$ or $0.5 \leq \alpha, \beta \leq 1$, $W(\alpha, \beta) = 0.4$,

otherwise; (3) Random geometric graphon: $W(\alpha, \beta) = f(\min(|\beta - \alpha|, 1 - |\beta - \alpha|))$, where
 $f(x) = e^{-\frac{x}{0.5-x}}$.

For the RL algorithms, we use the implementation of RLlib [36] (version 1.11.0, Apache-2.0 license). For PPO used to learn an optimal policy ensemble in block GMFC, we use a 64-dimensional linear layer to encode the observation and 2-layer MLPs with 256 hidden units per layer for both value network and actor network. For independent DQN and independent PPO, we use the default weight-sharing model with 64-dimensional embedding layers. We train the GMFC PPO agent for 1000 iterations, and other three MARL agents for 200 iterations. The specific hyper-parameters are listed in Table 3.

Table 3: RL Algorithm Settings

Algorithms	GMFC PPO	I-DQN	I-PPO	QMIX
Learning rate	0.0005	0.0005	0.0001	0.00005
Learning rate decay	True	True	True	False
Discount factor	0.95	0.95	0.95	0.95
Batch size	128	128	128	128
KL coefficient	0.2	-	0.2	-
KL target	0.01	-	0.01	-
Buffer size	-	2000	-	2000
Target network update frequency	-	2000	-	1000

6. Conclusion

In this work, we have proposed a discrete-time GMFC framework for MARL with nonuniform interactions and heterogeneous reward functions and transition functions across the agents on dense graphs. Theoretically, we have shown that under suitable assumptions, GMFC approximates MARL well with approximation error of order $\mathcal{O}(\frac{1}{\sqrt{N}})$. To reduce the dimension of GMFC, we have introduced block GMFC by discretizing the graphon index and shown that it also approximates MARL well. Empirical studies on several examples have verified the plausibility of the GMFC framework. For future research, we wish to explore more on how to extract the optimal policy of cooperative MARL without the simulator for population state distribution ensemble and to extend our framework to heterogeneous MARL on sparse graphs.

References

- [1] Adler, J. L., Blue, V. J., 2002. A cooperative multi-agent transportation management and route guidance system. *Transportation Research Part C: Emerging Technologies* 10 (5-6), 433–454.
- [2] Angiuli, A., Fouque, J.-P., Laurière, M., 2022. Unified reinforcement Q-learning for mean field game and control problems. *Mathematics of Control, Signals, and Systems*, 1–55.
- [3] Bayraktar, E., Chakraborty, S., Wu, R., 2020. Graphon mean field systems. arXiv preprint arXiv:2003.13180.

- 524 [4] Bet, G., Coppini, F., Nardi, F. R., 2020. Weakly interacting oscillators on dense random graphs.
525 arXiv preprint arXiv:2006.07670.
- 526 [5] Caines, P. E., Huang, M., 2019. Graphon mean field games and the GMFG equations: ε -Nash
527 equilibria. In: 2019 IEEE 58th conference on decision and control (CDC). IEEE, pp. 286–292.
- 528 [6] Cardaliaguet, P., Lehalle, C.-A., 2018. Mean field game of controls and an application to trade
529 crowding. *Mathematics and Financial Economics* 12 (3), 335–363.
- 530 [7] Carmona, R., Cooney, D. B., Graves, C. V., Lauriere, M., 2022. Stochastic graphon games: I.
531 the static case. *Mathematics of Operations Research* 47 (1), 750–778.
- 532 [8] Carmona, R., Delarue, F., 2013. Probabilistic analysis of mean-field games. *SIAM Journal on*
533 *Control and Optimization* 51 (4), 2705–2734.
- 534 [9] Carmona, R., Laurière, M., Tan, Z., 2019. Model-free mean-field reinforcement learning: mean-
535 field MDP and mean-field Q-learning. arXiv preprint arXiv:1910.12802.
- 536 [10] Chen, Y., Liu, J., Khossainov, B., 2021. Agent-level maximum entropy inverse reinforcement
537 learning for mean field games. arXiv preprint arXiv:2104.14654.
- 538 [11] Choi, J., Oh, S., Horowitz, R., 2009. Distributed learning and cooperative control for multi-
539 agent systems. *Automatica* 45 (12), 2802–2814.
- 540 [12] Cortes, J., Martinez, S., Karatas, T., Bullo, F., 2004. Coverage control for mobile sensing
541 networks. *IEEE Transactions on Robotics and Automation* 20 (2), 243–255.
- 542 [13] Cui, J., Liu, Y., Nallanathan, A., 2019. Multi-agent reinforcement learning-based resource
543 allocation for UAV networks. *IEEE Transactions on Wireless Communications* 19 (2), 729–
544 743.
- 545 [14] Cui, K., Koepl, H., 2021. Approximately solving mean field games via entropy-regularized deep
546 reinforcement learning. In: *International Conference on Artificial Intelligence and Statistics*.
547 PMLR, pp. 1909–1917.
- 548 [15] Cui, K., Koepl, H., 2022. Learning graphon mean field games and approximate Nash equilibria.
549 In: *International Conference on Learning Representations*.
- 550 [16] Cui, K., Tahir, A., Sinzger, M., Koepl, H., 2021. Discrete-time mean field control with envi-
551 ronment states. In: 2021 60th IEEE Conference on Decision and Control (CDC). IEEE, pp.
552 5239–5246.
- 553 [17] Delarue, F., 2017. Mean field games: A toy model on an Erdős-Renyi graph. *ESAIM: Proceed-*
554 *ings and Surveys* 60, 1–26.
- 555 [18] Elie, R., Perolat, J., Laurière, M., Geist, M., Pietquin, O., 2020. On the convergence of model
556 free learning in mean field games. In: *Proceedings of the AAAI Conference on Artificial Intel-*
557 *ligence*. Vol. 34. pp. 7143–7150.
- 558 [19] Esunge, J. N., Wu, J., 2014. Convergence of weighted empirical measures. *Stochastic Analysis*
559 *and Applications* 32 (5), 802–819.
- 560 [20] Flyvbjerg, H., Sneppen, K., Bak, P., 1993. Mean field theory for a simple model of evolution.
561 *Physical review letters* 71 (24), 4087.
- 562 [21] Gao, S., Caines, P. E., 2019. Graphon control of large-scale networks of linear systems. *IEEE*
563 *Transactions on Automatic Control* 65 (10), 4090–4105.
- 564 [22] Gao, S., Caines, P. E., 2019. Spectral representations of graphons in very large network systems
565 control. In: 2019 IEEE 58th conference on decision and Control (CDC). IEEE, pp. 5068–5075.
- 566 [23] Gu, H., Guo, X., Wei, X., Xu, R., 2021. Mean-field controls with Q-learning for cooperative
567 MARL: convergence and complexity analysis. *SIAM Journal on Mathematics of Data Science*
568 3 (4), 1168–1196.

- 569 [24] Gu, H., Guo, X., Wei, X., Xu, R., 2023. Dynamic programming principles for mean-field
570 controls with learning. *Operations Research*.
- 571 [25] Guo, X., Hu, A., Xu, R., Zhang, J., 2019. Learning mean-field games. *Advances in Neural
572 Information Processing Systems* 32.
- 573 [26] Hadikhanloo, S., Silva, F. J., 2019. Finite mean field games: fictitious play and convergence to
574 a first order continuous mean field game. *Journal de Mathématiques Pures et Appliquées* 132,
575 369–397.
- 576 [27] Huang, M., Malhamé, R. P., Caines, P. E., 2006. Large population stochastic dynamic games:
577 closed-loop McKean-Vlasov systems and the Nash certainty equivalence principle. *Communi-
578 cations in Information & Systems* 6 (3), 221–252.
- 579 [28] Hughes, E., Leibo, J. Z., Phillips, M., Tuyls, K., Dueñez-Guzman, E., García Castañeda, A.,
580 Dunning, I., Zhu, T., McKee, K., Koster, R., et al., 2018. Inequity aversion improves cooper-
581 ation in intertemporal social dilemmas. *Advances in neural information processing systems*
582 31.
- 583 [29] Lacker, D., Soret, A., 2022. A label-state formulation of stochastic graphon games and approx-
584 imate equilibria on large networks. *Mathematics of Operations Research*.
- 585 [30] Lasry, J.-M., Lions, P.-L., 2007. Mean field games. *Japanese Journal of Mathematics* 2 (1),
586 229–260.
- 587 [31] Lee, J. W., Park, J., Jangmin, O., Lee, J., Hong, E., 2007. A multiagent approach to Q -learning
588 for daily stock trading. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems
589 and Humans* 37 (6), 864–877.
- 590 [32] Lee, J. W., Zhang, B.-T., 2002. Stock trading system using reinforcement learning with cooper-
591 ative agents. In: *Proceedings of the Nineteenth International Conference on Machine Learning*.
592 pp. 451–458.
- 593 [33] Leibo, J. Z., Zambaldi, V., Lanctot, M., Marecki, J., Graepel, T., 2017. Multi-agent reinforce-
594 ment learning in sequential social dilemmas. *arXiv preprint arXiv:1702.03037*.
- 595 [34] Lerer, A., Peysakhovich, A., 2017. Maintaining cooperation in complex social dilemmas using
596 deep reinforcement learning. *arXiv preprint arXiv:1707.01068*.
- 597 [35] Li, Y., Wang, L., Yang, J., Wang, E., Wang, Z., Zhao, T., Zha, H., 2021. Permutation invariant
598 policy optimization for mean-field multi-agent reinforcement learning: A principled approach.
599 *arXiv preprint arXiv:2105.08268*.
- 600 [36] Liang, E., Liaw, R., Nishihara, R., Moritz, P., Fox, R., Goldberg, K., Gonzalez, J. E., Jor-
601 dan, M. I., Stoica, I., 2018. RLlib: Abstractions for distributed reinforcement learning. In:
602 *International Conference on Machine Learning (ICML)*.
- 603 [37] Lin, Y., Qu, G., Huang, L., Wierman, A., 2021. Multi-agent reinforcement learning in stochastic
604 networked systems. *Advances in Neural Information Processing Systems* 34, 7825–7837.
- 605 [38] Liu, B., Cai, Q., Yang, Z., Wang, Z., 2019. Neural proximal/trust region policy optimization
606 attains globally optimal policy. *arXiv preprint arXiv:1906.10306*.
- 607 [39] Lovász, L., 2012. Large networks and graph limits. Vol. 60. *American Mathematical Soc.*
- 608 [40] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller,
609 M., 2013. Playing Atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.
- 610 [41] Mondal, W., Aggarwal, V., Ukkusuri, S., 2022. On the near-optimality of local policies in large
611 cooperative multi-agent reinforcement learning. *Transactions on Machine Learning Research*.
- 612 [42] Mondal, W. U., Agarwal, M., Aggarwal, V., Ukkusuri, S. V., 2022. On the approximation of
613 cooperative heterogeneous multi-agent reinforcement learning (MARL) using mean field control
614 (MFC). *Journal of Machine Learning Research* 23 (129), 1–46.

- 615 [43] Mondal, W. U., Aggarwal, V., Ukkusuri, S. V., 2022. Can mean field control (MFC) approx-
616 imate cooperative multi agent reinforcement learning (MARL) with non-uniform interaction?
617 In: *Uncertainty in Artificial Intelligence*. PMLR, pp. 1371–1380.
- 618 [44] Pasztor, B., Bogunovic, I., Krause, A., 2021. Efficient model-based multi-agent mean-field
619 reinforcement learning. arXiv preprint arXiv:2107.04050.
- 620 [45] Qu, G., Wierman, A., Li, N., 2020. Scalable reinforcement learning of localized policies for
621 multi-agent networked systems. In: *Learning for Dynamics and Control*. PMLR, pp. 256–266.
- 622 [46] Rashid, T., Samvelyan, M., Schroeder, C., Farquhar, G., Foerster, J., Whiteson, S., 2018.
623 QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning. In:
624 *International Conference on Machine Learning*. PMLR, pp. 4295–4304.
- 625 [47] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O., 2017. Proximal policy opti-
626 mization algorithms. arXiv preprint arXiv:1707.06347.
- 627 [48] Subramanian, J., Mahajan, A., 2019. Reinforcement learning in stationary mean-field games.
628 In: *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent*
629 *Systems*. pp. 251–259.
- 630 [49] Vasal, D., Mishra, R., Vishwanath, S., 2021. Sequential decomposition of graphon mean field
631 games. In: *2021 American Control Conference (ACC)*. IEEE, pp. 730–736.
- 632 [50] W Axhausen, K., Horni, A., Nagel, K., 2016. *The multi-agent transport simulation MATSim*.
633 Ubiquity Press.
- 634 [51] Wainwright, M. J., 2019. *High-dimensional statistics: A non-asymptotic viewpoint*. Vol. 48.
635 Cambridge University Press.
- 636 [52] Yang, Y., Luo, R., Li, M., Zhou, M., Zhang, W., Wang, J., 2018. Mean field multi-agent
637 reinforcement learning. In: *International Conference on Machine Learning*. PMLR, pp. 5571–
638 5580.
- 639 [53] Yin, H., Mehta, P. G., Meyn, S. P., Shanbhag, U. V., 2010. Learning in mean-field oscillator
640 games. In: *49th IEEE Conference on Decision and Control (CDC)*. IEEE, pp. 3125–3132.

641 **Declaration of Competing Interest**

642 The authors declare that they have no known competing financial interests or personal
643 relationships that could have appeared to influence the work reported in this paper.

Journal Pre-proof