



**HAL**  
open science

# Quantitative propagation of chaos for mean field Markov decision process with common noise

Médéric Motte, Huyên Pham

► **To cite this version:**

Médéric Motte, Huyên Pham. Quantitative propagation of chaos for mean field Markov decision process with common noise. 2022. hal-03737655

**HAL Id: hal-03737655**

**<https://hal.science/hal-03737655>**

Preprint submitted on 25 Jul 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Quantitative propagation of chaos for mean field Markov decision process with common noise

Médéric MOTTE \*      Huyên PHAM †

July 25, 2022

## Abstract

We investigate propagation of chaos for mean field Markov Decision Process with common noise (CMKV-MDP), and when the optimization is performed over randomized open-loop controls on infinite horizon. We first state a rate of convergence of order  $M_N^\gamma$ , where  $M_N$  is the mean rate of convergence in Wasserstein distance of the empirical measure, and  $\gamma \in (0, 1]$  is an explicit constant, in the limit of the value functions of  $N$ -agent control problem with asymmetric open-loop controls, towards the value function of CMKV-MDP. Furthermore, we show how to explicitly construct  $(\epsilon + \mathcal{O}(M_N^\gamma))$ -optimal policies for the  $N$ -agent model from  $\epsilon$ -optimal policies for the CMKV-MDP. Our approach relies on sharp comparison between the Bellman operators in the  $N$ -agent problem and the CMKV-MDP, and fine coupling of empirical measures.

## 1 Introduction

We consider a social planner problem with  $N$  cooperative agents in a mean-field discrete time model with common noise over an infinite horizon. The controlled state process  $\mathbf{X} = (X^i)_{i \in \llbracket 1, N \rrbracket}$  of the  $N$ -agent model is given by the dynamical random system

$$\begin{cases} X_0^i &= x_0^i, \\ X_{t+1}^i &= F(X_t^i, \alpha_t^i, \frac{1}{N} \sum_{j=1}^N \delta_{(X_t^j, \alpha_t^j)}, \varepsilon_{t+1}^i, \varepsilon_{t+1}^0), \quad t \in \mathbb{N}. \end{cases} \quad (1.1)$$

Here,  $x_0^i$ ,  $i \in \llbracket 1, N \rrbracket$ , are the initial states valued in a compact Polish space  $\mathcal{X}$  with metric  $d$ ,  $(\varepsilon_t^i)_{i \in \llbracket 1, N \rrbracket, t \in \mathbb{N}^*}$  is a family of mutually i.i.d. random variables on some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , valued in some measurable space  $E$ , and representing idiosyncratic noises, while  $(\varepsilon_t^0)_{t \in \mathbb{N}^*}$  is another family of i.i.d. random variables valued in some measurable space  $E^0$ , and representing the common noise (independent of idiosyncratic noise). The control  $\alpha^i$  followed by agent  $i$ , is a process, valued in some compact Polish space  $A$  with metric  $d_A$ , and adapted with respect to the filtration  $(\mathcal{F}_t^N)_{t \in \mathbb{N}}$  generated by  $\varepsilon = ((\varepsilon_t^i)_{i \in \llbracket 1, N \rrbracket}, \varepsilon_t^0)_{t \in \mathbb{N}^*}$  and also completed with a family of

---

\*LPSM, Université Paris Cité [medericmotte at gmail.com](mailto:medericmotte@gmail.com)

†LPSM, Université Paris Cité, and CREST-ENSAE, [pham at lpsm.paris](mailto:pham@lpsm.paris) The author acknowledges support of the ANR 18-IDEX-001. This work was also partially supported by the Chair Finance & Sustainable Development / the FiME Lab (Institut Europlace de Finance)

mutually i.i.d. uniform random variables  $\mathbf{U} = (U_t^i)_{i \in \llbracket 1, N \rrbracket, t \in \mathbb{N}}$  that are used for randomization of the controls. The mean-field interaction between the agents is formalized via the state transition function  $F$  by the dependence upon the empirical measure of both state/action of all the other agents: here  $F$  is a measurable function from  $\mathcal{X} \times A \times \mathcal{P}(\mathcal{X} \times A) \times E \times E^0$  into  $\mathcal{X}$ , where  $\mathcal{P}(\mathcal{X} \times A)$  is the space of probability measures on the product space  $\mathcal{X} \times A$ .

The objective of the social planner is to maximize over the set  $\mathcal{A}$  of  $A^N$ -valued  $(\mathcal{F}_t^N)_{t \in \mathbb{N}}$ -adapted processes  $\alpha = (\alpha_t^i)_{i \in \llbracket 1, N \rrbracket, t \in \mathbb{N}}$  a criterion in the form

$$V_N^\alpha(\mathbf{x}_0) := \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{\infty} \beta^t f(X_t^i, \alpha_t^i, \frac{1}{N} \sum_{j=1}^N \delta_{(X_t^j, \alpha_t^j)}) \right],$$

where we set  $\mathbf{x}_0 = (x_0^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}^N$  for the initial state of the  $N$  agent system. Here  $\beta \in (0, 1)$  is a discount factor, and  $f$  is a bounded measurable real-valued function on  $\mathcal{X} \times A \times \mathcal{P}(\mathcal{X} \times A)$ . The value function for this optimization problem is defined on  $\mathcal{X}^N$  as

$$V_N(\mathbf{x}_0) := \sup_{\alpha \in \mathcal{A}} V_N^\alpha(\mathbf{x}_0), \quad (1.2)$$

and we notice that problem (1.1)-(1.2) is a standard Markov Decision Process (MDP) with state space  $\mathcal{X}^N$ , action space  $A^N$ , and (randomized) open-loop controls, and is the mathematical framework for reinforcement learning with multiple agents in interaction.

Let us now formulate the asymptotic mean-field problem when the number of agents  $N$  goes to infinity. This consists formally in replacing empirical distributions by theoretical ones in the dynamic system and gain functions. The controlled state process  $X$  of the representative agent is given by

$$\begin{cases} X_0 &= \xi_0, \\ X_{t+1} &= F(X_t^\alpha, \alpha_t, \mathbb{P}_{(X_t, \alpha_t)}^0, \varepsilon_{t+1}, \varepsilon_{t+1}^0), \quad t \in \mathbb{N}, \end{cases} \quad (1.3)$$

where we have renamed the uniform random sequence  $(U_t^1)_{t \in \mathbb{N}}$  and the noise  $(\varepsilon_t^1)_{t \in \mathbb{N}}$  by  $(U_t)_{t \in \mathbb{N}}$  and  $(\varepsilon_t)_{t \in \mathbb{N}}$ , and the initial state  $\xi_0$  is a  $\mathcal{G}$ -measurable random variable, with  $\mathcal{G}$  a  $\sigma$ -algebra independent of  $(U_t)_{t \in \mathbb{N}}$ ,  $(\varepsilon_t)_{t \in \mathbb{N}}$ ,  $(\varepsilon_t^0)_{t \in \mathbb{N}^*}$ , with distribution law  $\mu_0 \in \mathcal{P}(\mathcal{X})$  (the set of probability measures on  $\mathcal{X}$ ). The control process  $\alpha$  is an  $A$ -valued process, adapted with respect to the filtration generated by  $\mathcal{G}$ ,  $(U_t)_{t \in \mathbb{N}}$ ,  $(\varepsilon_t)_{t \in \mathbb{N}}$ ,  $(\varepsilon_t^0)_{t \in \mathbb{N}^*}$ , denoted by  $\alpha \in \mathcal{A}$ . Here  $\mathbb{P}^0$  and  $\mathbb{E}^0$  represent the conditional probability and expectation knowing the common noise  $\varepsilon^0$ , and then, given a random variable  $Y$ , we denote by  $\mathbb{P}_Y^0$  or  $\mathcal{L}^0(Y)$  its conditional law knowing  $\varepsilon^0$ . The McKean-Vlasov (or mean-field) control problem consists in maximizing over randomized open-loop controls  $\alpha$  in  $\mathcal{A}$  the gain functional

$$V^\alpha(\xi_0) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \beta^t f(X_t, \alpha_t, \mathbb{P}_{(X_t, \alpha_t)}^0) \right].$$

The value function to this optimization problem is defined on  $\mathcal{P}(\mathcal{X})$  by

$$V(\xi_0) := \sup_{\alpha \in \mathcal{A}} V^\alpha(\xi_0), \quad (1.4)$$

and we recall from [16] that  $V$  depends on  $\xi_0$  only through its distribution (invariance in law), and we denote by misuse of notation:  $V(\mu_0) = V(\xi_0)$ . Problem (1.3)-(1.4) is called mean-field

Markov Decision Process with common noise, or conditional McKean-Vlasov Markov Decision Process (CMKV-MDP in short), with the peculiarity compared to standard MDP coming from the dependence of the state transition on the conditional distribution of the state/action. In view of propagation of chaos for particle systems usually derived for mean-field diffusion process (see [18]), it is expected that CMKV-MDP provides a mean-field approximation of the  $N$ -agent MDP model.

While the literature on mean-field control in continuous time, in particular the optimal control of McKean-Vlasov equations, is quite important, see the monograph [5] for an overview and related references, there are rather few papers devoted to the discrete time framework. One of the first works is [11] which studies the convergence of large interacting population process to a simple mean-field model when the state space is finite. The paper [17] studies a discrete-time McKean-Vlasov control problem with feedback controls on finite horizon, and derive the corresponding dynamic programming equation which is explicitly solved in the linear quadratic case. In [6], the authors consider mean-field control on infinite horizon with common noise with a discussion about connections between closed-loop and open-loop policies, and propose  $Q$ -learning algorithms. Our companion paper [16] deals with open-loop control and highlights the role of randomized controls with respect to standard Markov Decision Process (MDP). The value function is characterized as a fixed point Bellman equation defined on the space of probability measures, and existence of  $\epsilon$ -optimal randomized feedback controls is proved. The recent paper [1] studies mean-field control with deterministic closed-loop policies through the lens of MDP theory, and discusses the existence of optimal policies for the limiting mean-field problem as well as for the  $N$ -agent problem.

**Main contributions.** In this paper, we establish a quantitative propagation of result for the  $N$ -agent MDP towards the CMKV-MDP. Our contributions are twofold:

1. We show in Theorem 2.1 an explicit rate of convergence of the value functions under some assumptions to be precised later: there exists some positive constant  $C$  (depending on the data of the problem) such that for all  $\mathbf{x} = (x^i)_{i \in [1, N]} \in \mathcal{X}^N$ ,

$$\left| V_N(\mathbf{x}) - V\left(\frac{1}{N} \sum_{i=1}^N \delta_{x^i}\right) \right| \leq CM_N^\gamma,$$

where  $M_N$  is the mean rate of convergence in Wasserstein distance of the empirical measure (see [9]), and  $\gamma \in (0, 1]$  is an explicit constant depending on  $\beta$  and  $F$ .

2. We prove that any  $\epsilon$ -optimal randomized feedback policy for the CMKV-MDP (including the case  $\epsilon = 0$ , i.e., optimal randomized feedback policy whose existence is shown) yields either an approximate optimal feedback control or an approximate randomized feedback control for the  $N$ -agent MDP problem, in a constructive sense to be precised later with an explicit rate of convergence, see Theorems 2.2 and 2.3.

While the first statement for convergence of value function is important in theory, the second statement is particularly interesting in practice (but often less studied in the literature) since it means that if the McKean-Vlasov MDP is simpler to solve than the  $N$ -agent MDP (some examples and applications to targeted advertising are developed in the PhD thesis [15]), then one can compute an almost optimal randomized feedback policy for the McKean-Vlasov MDP, and then use it in the  $N$ -agent MDP: this will guaranty us to have an almost optimal control.

*Related literature.* The convergence of the  $N$ -individual problem to the limiting mean-field control problem has been first rigorously proved in [14] by tightness and martingale arguments for continuous-time controlled McKean-Vlasov equations. This result has been extended in [8] to the common noise case and when there is interaction via the joint distribution of the state and control. The paper [10] proved by viscosity solutions method via the characterization of the Hamilton-Jacobi-Bellman equation the convergence of the value function towards the  $N$ -agent problem to the value function of the mean-field control problem in the common noise case but without idiosyncratic noise. Rate of convergence of order  $1/N$  has been stated in [12] by Backward Stochastic Differential Equations techniques but under the strong condition that there exists a smooth solution to the Master Bellman equation. The recent paper [4] removed this regularity assumption on the value function, and obtained an algebraic rate of convergence of order  $N^{-\gamma}$  for some constant  $\gamma \in (0, 1]$ . We mention also in the continuous-time framework the paper [7] which derived a rate of convergence of order  $N^{-1/2}$  when the state space is finite.

The convergence of the value function in the  $N$ -agent problem in a discrete-time mean field framework has been studied in our companion paper [16]. However, it was assumed there that each agent used the same open-loop policy, applied to her own idiosyncratic noise and the common noise. In particular, agent's controls cannot depend upon other agent's idiosyncratic noises, and they have symmetric (or exchangeable) behaviours. This restriction was crucial for using propagation of chaos argument relying on a pathwise comparison between the state and control processes in the  $N$ -individual model and the McKean-Vlasov MDPs.

In this paper, we consider that the control of each agent can also depend upon the idiosyncratic noises of all the population, and that they can do so in a completely asymmetric way (i.e. each agent can use a different open-loop policy). This additional flexibility and generality in the definition of controls prevents us from coupling controls between the  $N$ -agent and the McKean-Vlasov MDPs in a one-to-one fashion as in [16]. In order to overcome this difficulty, we adopt quite different arguments by coupling the Bellman operators instead of the state/control process of the  $N$ -agent and CMKV MDPs. More precisely, the strategy of the proof is the following:

*Idea of the proof.*

- (i) We first derive the Bellman equation for the  $N$ -agent MDP, with arguments similar to [16], i.e. we prove that  $\mathcal{T}_N V_N = V_N$ , where  $\mathcal{T}_N$  is the operator defined by

$$\mathcal{T}_N W(\mathbf{x}) := \sup_{\mathbf{a} \in A^N} \mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X}^N,$$

with

$$\mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f(x^i, a^i, \frac{1}{N} \sum_{j=1}^N \delta_{(x^j, a^j)}) + \beta \mathbb{E}[W((F(x^i, a^i, \frac{1}{N} \sum_{j=1}^N \delta_{(x^j, a^j)}, \varepsilon_1^i, \varepsilon_1^0)_{i \in \llbracket 1, N \rrbracket}))],$$

for  $\mathbf{x} = (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}^N$ ,  $\mathbf{a} = (a^i)_{i \in \llbracket 1, N \rrbracket} \in A^N$ . This property is obtained by seeing the  $N$ -agent MDP as a standard MDP on  $\mathcal{X}^N$  with actions space  $A^N$ .

- (ii) Then, we observe that the operators  $\mathbb{T}^{\mathbf{a}}$  of the McKean-Vlasov MDP, derived in [16], are, formally, the limits of  $\mathbb{T}_N^{\mathbf{a}}$  when  $N \rightarrow \infty$ , for  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1], A)$  and  $\mathbf{a} \in A^N$  well coupled. Inspired by this formal observation, we “compare”  $\mathbb{T}_N^{\mathbf{a}}$  to  $\mathbb{T}^{\mathbf{a}}$  and prove that they are indeed “close” in some sense, for  $N$  large. A key point is that  $\mathbb{T}_N^{\mathbf{a}}$  is defined on  $L_m^\infty(\mathcal{X}^N)$  (the set

of bounded measurable functions on  $\mathcal{X}^N$ , valued in  $\mathbb{R}$ ) while  $\mathbb{T}^a$  is defined on  $L_m^\infty(\mathcal{P}(\mathcal{X}))$  (the set of bounded measurable functions on  $\mathcal{P}(\mathcal{X})$ , valued in  $\mathbb{R}$ ). To compare both type of objects, we introduce a canonical way to associate to a function  $W \in L_m^\infty(\mathcal{P}(\mathcal{X}))$  the function  $\widetilde{W} \in L_m^\infty(\mathcal{X}^N)$  by setting  $\widetilde{W}(\mathbf{x}) = W\left(\frac{1}{N} \sum_{i=1}^N \delta_{x^i}\right)$ .

- (iii) Once the proximity between  $\mathbb{T}_N^{\mathbf{a}}$  and  $\mathbb{T}^a$  is established in a general sense, we prove the proximity of the value functions  $V_N$  and  $V$  by seeing them as the unique fixed points of the Bellman operators  $\mathcal{T}_N = \sup_{\mathbf{a} \in A^N} \mathbb{T}_N^{\mathbf{a}}$  and  $\mathcal{T} = \sup_{\mathbf{a} \in L^0(\mathcal{X} \times [0,1], A)} \mathbb{T}^{\mathbf{a}}$ , following the intuition that if two contracting operators are close, their unique fixed points should also be close.
- (iv) Finally, we provide two procedures to build  $\mathcal{O}(\epsilon + M_N^\gamma)$ -optimal policies for the  $N$ -agent MDP from an  $\epsilon$ -optimal stationary randomized feedback policy for the McKean-Vlasov MDP. The idea is to view, for each MDP, any  $\epsilon$ -optimal policy as a policy satisfying the verification theorem, which is a property only linked to the Bellman operator, again following the intuition that if two Bellman operators are close, the policies satisfying their verification results should also be close.

**Outline of the paper.** The rest of the paper is organized as follows. We state the assumptions and the main results in Section 2, while Section 3 is devoted to their proofs. Finally, we give in Appendix A the proof of existence for optimal randomized feedback policy, and put in Appendix B some results about the Bellman operator for the  $N$ -agent MDP problem that are needed in the proof of our convergence results.

## 2 Main results

### 2.1 Notations and assumptions

The product space  $\mathcal{X} \times A$  is equipped with the metric  $\mathbf{d}((x, a), (x', a')) = d(x, x') + d_A(a, a')$ ,  $x, x' \in \mathcal{X}$ ,  $a, a' \in A$ . Likewise, we shall endow  $\mathcal{X}^N$  with the metric  $\mathbf{d}_N(\mathbf{x}, \mathbf{x}') = \frac{1}{N} \sum_{i=1}^N d(x^i, x'^i)$  for  $\mathbf{x} = (x^i)_{i \in [1, N]}$ ,  $\mathbf{x}' = (x'^i)_{i \in [1, N]} \in \mathcal{X}^N$ ,  $A^N$  with the metric  $\mathbf{d}_{A, N}(\mathbf{a}, \mathbf{a}') = \frac{1}{N} \sum_{i=1}^N d_A(a^i, a'^i)$  for  $\mathbf{a} = (a^i)_{i \in [1, N]}$ ,  $\mathbf{a}' = (a'^i)_{i \in [1, N]} \in A^N$ , and  $(\mathcal{X} \times A)^N$  with the metric  $\mathbf{d}_N((\mathbf{x}, \mathbf{a}), (\mathbf{x}', \mathbf{a}')) = \frac{1}{N} \sum_{i=1}^N \mathbf{d}((x^i, a^i), (x'^i, a'^i))$  for  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^N$  and  $\mathbf{a}, \mathbf{a}' \in A^N$ . When  $(\mathcal{Y}, d)$  is a compact metric space, the set  $\mathcal{P}(\mathcal{Y})$  of probability measures on  $\mathcal{Y}$  is equipped with the Wasserstein distance

$$\mathcal{W}_d(\mu, \mu') := \inf \left\{ \int_{\mathcal{Y}^2} d(y, y') \mu(dy, dy') : \mu \in \mathbf{\Pi}(\mu, \mu') \right\},$$

where  $\mathbf{\Pi}(\mu, \mu')$  is the set of (coupling) probability measures on  $\mathcal{Y} \times \mathcal{Y}$  with marginals  $\mu$  and  $\mu'$ , and we recall the dual Kantorovich-Rubinstein representation

$$\mathcal{W}_d(\mu, \mu') = \sup_{\phi \in Lip_1} \int_{\mathcal{Y}} \phi(y) (\mu - \mu')(dy), \quad (2.1)$$

where  $Lip_1$  is the set of Lipschitz functions on  $\mathcal{Y}$  with Lipschitz constant bounded by 1.

Given  $\mathbf{x} = (x^i)_{i \in [1, N]} \in \mathcal{X}^N$ , and  $\mathbf{a} = (a^i)_{i \in [1, N]} \in A^N$ , we denote by

$$\mu_N[\mathbf{x}] := \frac{1}{N} \sum_{i=1}^N \delta_{x^i} \in \mathcal{P}(\mathcal{X}), \quad \mu_N[\mathbf{x}, \mathbf{a}] := \frac{1}{N} \sum_{i=1}^N \delta_{(x^i, a^i)} \in \mathcal{P}(\mathcal{X} \times A),$$

and we recall that

$$\mathcal{W}_d(\mu_N[\mathbf{x}, \mathbf{a}], \mu_N[\mathbf{x}', \mathbf{a}']) \leq \mathbf{d}_N((\mathbf{x}, \mathbf{a}), (\mathbf{x}', \mathbf{a}')). \quad (2.2)$$

Given a random variable  $Y$  on  $(\Omega, \mathcal{F}, \mathbb{P})$ , we denote by  $\mathbb{P}_Y$  or  $\mathcal{L}(Y)$  its distribution law.

We make the following standing assumptions on the state transition function  $F$  and on the running reward function  $f$ .

**(HF<sub>lip</sub>)** There exists  $K_F > 0$ , such that for all  $a, a' \in A$ ,  $e^0 \in E^0$ ,  $x, x' \in \mathcal{X}$ ,  $\mu, \mu' \in \mathcal{P}(\mathcal{X} \times A)$ ,

$$\mathbb{E}[d(F(x, a, \mu, \varepsilon_1^1, e^0), F(x', a', \mu', \varepsilon_1^1, e^0))] \leq K_F(\mathbf{d}((x, a), (x', a')) + \mathcal{W}_d(\mu, \mu')).$$

**(Hf<sub>lip</sub>)** There exists  $K_f > 0$ , such that for all  $x, x' \in \mathcal{X}$ ,  $a, a' \in A$ ,  $\mu, \mu' \in \mathcal{P}(\mathcal{X} \times A)$ ,

$$|f(x, a, \mu) - f(x', a', \mu')| \leq K_f(\mathbf{d}((x, a), (x', a')) + \mathcal{W}_d(\mu, \mu')).$$

**Remark 2.1** We stress the importance of making the regularity assumptions for  $F$  in *expectation* only. When  $\mathcal{X}$  is finite,  $F$  cannot be, strictly speaking, Lipschitz (or even continuous) unless it is constant w.r.t. its mean-field argument ( $\mu$  and  $\mu'$  in **(HF<sub>lip</sub>)**). However,  $F$  can be Lipschitz *in expectation*, e.g. once integrated w.r.t. the idiosyncratic noise.

Under Assumption **(HF<sub>lip</sub>)**, we define the constant

$$\gamma := \min \left[ 1, \frac{|\ln \beta|}{\ln(2K_F)_+} \right] \in (0, 1].$$

In the sequel, we denote by  $\Delta_{\mathcal{X}}$  (resp.  $\Delta_A$  and  $\Delta_{\mathcal{X} \times A}$ ) the diameter of the compact metric space  $\mathcal{X}$  (resp.  $A$  and  $\mathcal{X} \times A$ ), and define

$$M_N := \sup_{\mu \in \mathcal{P}(\mathcal{X} \times A)} \mathbb{E}[\mathcal{W}_d(\mu_N, \mu)], \quad (2.3)$$

where  $\mu_N$  is the empirical measure  $\mu_N = \frac{1}{N} \sum_{n=1}^N \delta_{Y_n}$ ,  $(Y_n)_{1 \leq n \leq N}$  are i.i.d. random variables with law  $\mu$ . It is known that  $M_N \xrightarrow{N \rightarrow \infty} 0$ , and we recall from [9], and [3] some results about non asymptotic bounds for the mean rate of convergence in Wasserstein distance of the empirical measure.

- If  $\mathcal{X} \times A \subset \mathbb{R}^d$  for some  $d \in \mathbb{N}^*$ , then:  $M_N = \mathcal{O}(N^{-\frac{1}{2}})$  for  $d = 1$ ,  $M_N = \mathcal{O}(N^{-\frac{1}{2}} \log(1 + N))$  for  $d = 2$ , and  $M_N = \mathcal{O}(N^{-\frac{1}{d}})$  for  $d \geq 3$ .
- If for all  $\delta > 0$ , the smallest number of balls with radius  $\delta$  covering the compact metric set  $\mathcal{X} \times A$  with diameter  $\Delta_{\mathcal{X} \times A}$  is smaller than  $\mathcal{O}\left(\left(\frac{\Delta_{\mathcal{X} \times A}}{\delta}\right)^\theta\right)$  for  $\theta > 2$ , then  $M_N = \mathcal{O}(N^{-1/\theta})$ .

In the sequel  $C$  will denote a generic constant that depends only on the data of the problem, namely  $\Delta_{\mathcal{X}}$ ,  $\Delta_{\mathcal{X} \times A}$ ,  $\beta$ ,  $K_F$  and  $K_f$ .

## 2.2 Convergence of value functions

Our first main result is to quantify the rate of convergence of the value function of the  $N$ -agent MDP towards the value function of the CMKV-MDP.

**Theorem 2.1** *There exists some positive constant  $C$  such that for all  $\mathbf{x} = (x^i)_{i \in [1, N]} \in \mathcal{X}^N$ ,*

$$\left| V_N(\mathbf{x}) - V(\mu_N[\mathbf{x}]) \right| \leq CM_N^\gamma.$$

## 2.3 Approximate optimal policies

Our next results are to show how to obtain approximate optimal control for the  $N$  agent MDP from  $\varepsilon$ -optimal control for CKMV-MDP, and to quantify the accuracy of this approximation.

First, let us recall from [16] the construction of  $\varepsilon$ -optimal control for CKMV-MDP. The value function  $V$  is characterized as the unique fixed point in  $L_m^\infty(\mathcal{P}(\mathcal{X}))$ , the set of bounded measurable real-valued functions on  $\mathcal{P}(\mathcal{X})$ , of the Bellman equation  $V = \mathcal{T}V$ , where  $\mathcal{T}$  is the Bellman operator defined on  $L_m^\infty(\mathcal{P}(\mathcal{X}))$  by

$$\begin{aligned} \mathcal{T}W(\mu) &:= \sup_{\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)} \mathbb{T}^{\mathbf{a}}W(\mu), \\ \text{with } \mathbb{T}^{\mathbf{a}}W(\mu) &:= \mathbb{E} \left[ f(\xi, \mathbf{a}(\xi, U), \mathcal{L}(\xi, \mathbf{a}(\xi, U))) + \beta W(\mathbb{P}_{F(\xi, \mathbf{a}(\xi, U), \mathcal{L}(\xi, \mathbf{a}(\xi, U)), \varepsilon_1, \varepsilon_1^0)}^0) \right], \end{aligned} \quad (2.4)$$

for any  $(\xi, U) \sim \mu \otimes \mathcal{U}([0, 1])$  (it is clear that the right-hand side in (2.4) does not depend on the choice of such  $(\xi, U)$ ), where  $L^0(\mathcal{X} \times [0, 1]; A)$  is the set of measurable functions from  $\mathcal{X} \times [0, 1]$  into  $A$ . Then, for all  $\varepsilon > 0$ , there exists a randomized feedback policy  $\mathbf{a}_\varepsilon$ , i.e. a measurable function from  $\mathcal{P}(\mathcal{X}) \times \mathcal{X} \times [0, 1]$  into  $A$ , denoted by  $\mathbf{a}_\varepsilon \in L^0(\mathcal{P}(\mathcal{X}) \times \mathcal{X} \times [0, 1]; A)$ , such that for all  $\mu \in \mathcal{P}(\mathcal{X})$ :

$$V(\mu) - \varepsilon \leq \mathbb{T}^{\mathbf{a}_\varepsilon(\mu, \cdot)}V(\mu),$$

and we say that  $\mathbf{a}_\varepsilon$  is an  $\varepsilon$ -optimal randomized feedback policy for CMKV-MDP. By considering the randomized feedback control  $\alpha^\varepsilon \in \mathcal{A}$  defined by

$$\alpha_t^\varepsilon = \mathbf{a}_\varepsilon(\mathbb{P}_{X_t}^0, X_t, U_t), \quad t \in \mathbb{N}, \quad (2.5)$$

where  $(U_t)_{t \in \mathbb{N}}$  is an i.i.d. sequence of random variables,  $U_t \sim \mathcal{U}([0, 1])$ , independent of  $\xi_0 \sim \mu_0$ , and  $\varepsilon$ , this yields an  $O(\varepsilon)$ -optimal control for  $V(\mu_0)$ , namely

$$V(\mu_0) - \frac{\varepsilon}{1 - \beta} \leq V^{\alpha^\varepsilon}(\xi_0).$$

Actually, we can even take  $\varepsilon = 0$ , i.e., get optimal randomized feedback control. The proof for the existence of an optimal randomized feedback policy is inspired by the paper [6], which states the existence of an optimal policy in a closely related model, and is reported in Appendix A.

We now provide two procedures to construct approximate optimal control for the  $N$ -agent MDP from an  $\varepsilon$ -optimal randomized feedback policy for CMKV-MDP. The first procedure gives a general approach for getting approximate feedback control for the  $N$ -agent MDP.



**Theorem 2.2** *Let  $\mathbf{a}_\epsilon$  be an  $\epsilon$ -optimal randomized feedback policy for CMKV-MDP. Then, there exists a measurable function  $\boldsymbol{\pi}^{\mathbf{a}_\epsilon, N}$  from  $\mathcal{X}^N$  into  $A^N$ , called feedback policy for the  $N$ -agent MDP, such that*

$$\boldsymbol{\pi}^{\mathbf{a}_\epsilon, N}(\mathbf{x}) \in \underset{\mathbf{a} \in A^N}{\operatorname{argmin}} \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]), \quad \mathbf{x} \in \mathcal{X}^N, \quad (2.6)$$

with  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ . This yields a feedback control  $\boldsymbol{\alpha}^{\epsilon, N} \in \mathcal{A}$  defined by

$$\boldsymbol{\alpha}_t^{\epsilon, N} = \boldsymbol{\pi}^{\mathbf{a}_\epsilon, N}(\mathbf{X}_t), \quad t \in \mathbb{N},$$

which is  $O(\epsilon + M_N^\gamma)$ -optimal control for  $V_N(\mathbf{x}_0)$ , namely:

$$V_N(\mathbf{x}_0) - C[\epsilon + M_N^\gamma] \leq V_N^{\boldsymbol{\alpha}^{\epsilon, N}}(\mathbf{x}_0).$$

Theorem 2.2 provides a generic way to obtain a  $O(\epsilon + M_N^\gamma)$ -optimal feedback policy for the  $N$ -agent MDP from an  $\epsilon$ -optimal randomized feedback policy  $\mathbf{a}_\epsilon$  for CMKV-MDP, simply by sending actions  $\mathbf{a} = (a^i)_{i \in \llbracket 1, N \rrbracket}$  to the population so that, once in state  $\mathbf{x}$ , the state-action pair  $(\mathbf{x}, \mathbf{a})$  is empirically distributed as closely as possible to  $\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U))$ . However, the computation of this argmin in (2.6) can be difficult in practice.

We propose a second approach which provides a more practical derivation of an approximate optimal control for the  $N$ -agent MDP. It will use randomized feedback policy for the  $N$ -agent model, defined as a measurable function from  $\mathcal{X}^N \times [0, 1]^N$  into  $A^N$ .

**Theorem 2.3** *Let  $\mathbf{a}_\epsilon$  be an  $\epsilon$ -optimal randomized feedback policy for CMKV-MDP, assumed to satisfy the regularity condition*

$$\mathbb{E}[d_A(\mathbf{a}_\epsilon(\mu, x, U), \mathbf{a}_\epsilon(\mu, x', U))] \leq Kd(x, x'), \quad \forall x, x' \in \mathcal{X}, \mu \in \mathcal{P}(\mathcal{X}), \quad (2.7)$$

(here  $U \sim \mathcal{U}([0, 1])$ ) for some positive constant  $K$ . Consider the randomized feedback policy in the  $N$ -agent model defined by

$$\boldsymbol{\pi}_r^{\mathbf{a}_\epsilon, N}(\mathbf{x}, \mathbf{u}) := (\mathbf{a}_\epsilon(\mu_N[\mathbf{x}], x^i, u^i))_{i \in \llbracket 1, N \rrbracket},$$

for  $\mathbf{x} = (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}^N$ ,  $\mathbf{u} = (u^i)_{i \in \llbracket 1, N \rrbracket} \in [0, 1]^N$ . Then, the randomized feedback control  $\boldsymbol{\alpha}^{r, \epsilon, N} \in \mathcal{A}$  defined as

$$\boldsymbol{\alpha}_t^{r, \epsilon, N} = \boldsymbol{\pi}_r^{\mathbf{a}_\epsilon, N}(\mathbf{X}_t, \mathbf{U}_t), \quad t \in \mathbb{N},$$

where  $\{\mathbf{U}_t = (U_t^i)_{i \in \llbracket 1, N \rrbracket}, t \in \mathbb{N}\}$  is a family of mutually i.i.d. uniform random variables on  $[0, 1]$ , independent of  $\mathcal{G}$ ,  $\boldsymbol{\varepsilon} = ((\varepsilon_t^i)_{i \in \llbracket 1, N \rrbracket}, \varepsilon_t^0)_{t \in \mathbb{N}^*}$ , is an  $O(\epsilon + M_N^\gamma)$ -optimal control for  $V_N(\mathbf{x}_0)$ , namely:

$$V_N(\mathbf{x}_0) - C(1 + K)(\epsilon + M_N^\gamma) \leq V_N^{\boldsymbol{\alpha}^{r, \epsilon, N}}(\mathbf{x}_0).$$

Theorem 2.3 provides a simple and natural procedure to get an approximate policy for the  $N$ -agent MDP: it corresponds to using an  $\epsilon$ -optimal randomized feedback policy  $\mathbf{a}_\epsilon$  of the CMKV-MDP, but instead of inputting the theoretical state distribution of the McKean-Vlasov MDP in its mean-field argument, we input the empirical state distribution of the  $N$ -agent MDP, and instead of inputting the McKean-Vlasov state in its state argument, we input the  $N$ -agent individual states, and moreover, we use a randomization by tossing a coin at any time and for any agent. Notice that the validity of this procedure requires the Lipschitz condition (2.7), which always holds true when the state space  $\mathcal{X}$  is finite. Indeed, in this case, the metric on  $\mathcal{X}$  is the discrete distance  $d(x, x') = 1_{x \neq x'}$ , and (2.7) is clearly satisfied with  $K = \Delta_A$ .

### 3 Proof of main results

This section is devoted to the proofs of Theorems 2.1, 2.2, and 2.3 about rate of convergence in the propagation of chaos between the  $N$ -agent MDP and the limiting conditional McKean-Vlasov MDP. Our approach relies on the Bellman operators of each MDP. By proving their proximity (in a sense to be precised), we will be able to prove on the one hand the proximity of their unique fixed points, hence the convergence of the value functions, and on the other hand that almost optimal randomized feedback policies are directly related to the Bellman operators via the verification result, which will give the convergence of the approximate controls.

#### 3.1 Comparing the Bellman operators

We first introduce the following useful measurable optimal permutation for the coupling of empirical measures.

**Definition 3.1 (Measurable optimal permutation)** *Let  $(\mathcal{Y}, d)$  be a metric space. There exists a measurable map  $\sigma : (\mathbf{y}, \mathbf{y}') \in (\mathcal{Y}^N)^2 \rightarrow \sigma^{\mathbf{y}, \mathbf{y}'} \in \mathfrak{S}_N$  (where  $\mathfrak{S}_N$  denotes the set of permutations on  $\llbracket 1, N \rrbracket$ ) such that for all  $(\mathbf{y}, \mathbf{y}') \in (\mathcal{Y}^N)^2$ , we have*

$$\mathcal{W}_d(\mu_N[\mathbf{y}], \mu_N[\mathbf{y}']) = \mathbf{d}_N(\mathbf{y}, \mathbf{y}'_{\sigma^{\mathbf{y}, \mathbf{y}'}}), \quad (3.1)$$

where we set  $\mathbf{y}'_{\sigma^{\mathbf{y}, \mathbf{y}'}} = (y'^{\sigma_i^{\mathbf{y}, \mathbf{y}'}})_{i \in \llbracket 1, N \rrbracket}$  for  $\mathbf{y}' = (y'^i)_{i \in \llbracket 1, N \rrbracket}$ .

**Proof.** It is a well known result (see [19]) that, given  $(\mathbf{y}, \mathbf{y}') \in (\mathcal{Y}^N)^2$ , there exists a permutation  $\sigma^{\mathbf{y}, \mathbf{y}'} \in \mathfrak{S}_N$  realizing an optimal coupling between  $\mu_N[\mathbf{y}], \mu_N[\mathbf{y}'] \in \mathcal{P}(\mathcal{Y})$ , i.e., s.t. (3.1) holds. Let us check that this optimal permutation can be represented as a measurable function of  $(\mathbf{y}, \mathbf{y}') \in (\mathcal{Y}^N)^2$ . Let  $n \in \llbracket 1, N! \rrbracket \mapsto \sigma^n \in \mathfrak{S}_N$  be some bijection. Notice that the function

$$\mathbf{y}, \mathbf{y}' \in \mathcal{Y}^N \mapsto (\mathbf{d}_N(\mathbf{y}, \mathbf{y}'_{\sigma^n}))_{n \in \llbracket 1, N! \rrbracket} \in \mathbb{R}^{N!}$$

is continuous, hence measurable. Furthermore, it is clear that the function

$$\mathbf{z} \in \mathbb{R}^{N!} \mapsto \min_{n \in N!} [\operatorname{argmin} z^n]$$

is measurable. Denoting by

$$n_{\min}(\mathbf{y}, \mathbf{y}') := \min_{n \in N!} [\operatorname{argmin} \mathbf{d}_N(\mathbf{y}, \mathbf{y}'_{\sigma^n})],$$

it follows that the function  $\mathbf{y}, \mathbf{y}' \in \mathcal{X}^N \mapsto \sigma^{\mathbf{y}, \mathbf{y}'} = \sigma^{n_{\min}(\mathbf{y}, \mathbf{y}')}$  is a measurable representation of the optimal permutation.  $\square$

We now study the ‘‘proximity’’ between the Bellman operator of the CMKV-MDP given in (2.4), and the Bellman operator of the  $N$ -agent problem, viewed as a MDP with state space  $\mathcal{X}^N$ , action space  $A^N$ , noise sequence  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_t)_{t \in \mathbb{N}^*}$  with  $\boldsymbol{\varepsilon}_t := ((\varepsilon_t^i)_{i \in \llbracket 1, N \rrbracket}, \varepsilon_t^0)$  valued in  $E^N \times E^0$ , state transition function

$$\mathbf{F}(\mathbf{x}, \mathbf{a}, \mathbf{e}) := \left( F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], e^i, e^0) \right)_{i \in \llbracket 1, N \rrbracket}, \quad \mathbf{e} = ((e^i)_{i \in \llbracket 1, N \rrbracket}, e^0) \in E^N \times E^0,$$

and reward function

$$\mathbf{f}(\mathbf{x}, \mathbf{a}) = \frac{1}{N} \sum_{i=1}^N f(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}]), \quad \mathbf{x} = (x^i)_{i \in [1, N]}, \quad \mathbf{a} = (a^i)_{i \in [1, N]}.$$

Denoting by  $L_m^\infty(\mathcal{X}^N)$  the subset of measurable functions in  $L^\infty(\mathcal{X}^N)$  (the set of bounded real-valued functions on  $\mathcal{X}^N$ ), the Bellman “operator”  $\mathcal{T}_N : L_m^\infty(\mathcal{X}^N) \rightarrow L^\infty(\mathcal{X}^N)$  of the  $N$ -agent MDP is defined for any  $W \in L_m^\infty(\mathcal{X}^N)$  by:

$$\mathcal{T}_N W(\mathbf{x}) := \sup_{\mathbf{a} \in A^N} \mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X}^N,$$

where

$$\mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}) := \mathbf{f}(\mathbf{x}, \mathbf{a}) + \beta \mathbb{E}[W(\mathbf{F}(\mathbf{x}, \mathbf{a}, \varepsilon_1))], \quad \mathbf{x} \in \mathcal{X}^N, \quad \mathbf{a} \in A^N.$$

The characterization of the value function  $V_N$  and optimal controls for the  $N$ -agent MDP via the Bellman operator  $\mathcal{T}_N$  is stated in Appendix B.

We aim to quantify how “close”  $\mathbb{T}_N^{\mathbf{a}}$  and  $\mathbb{T}^{\mathbf{a}}$  are when  $\mathbf{a}$  and  $\mathbf{a}$  are close in a sense to be precised. Notice that the  $N$ -agent operator  $\mathbb{T}_N^{\mathbf{a}}$  is defined on  $L_m^\infty(\mathcal{X}^N)$  while the McKean-Vlasov operator  $\mathbb{T}^{\mathbf{a}}$  is defined on  $L_m^\infty(\mathcal{P}(\mathcal{X}))$ . There is however a natural way to compare them by means of an “unlifting” procedure. To any function  $W \in L_m^\infty(\mathcal{P}(\mathcal{X}))$ , we associate the unlifted function  $\widetilde{W} \in L_m^\infty(\mathcal{X}^N)$  defined by

$$\widetilde{W}(\mathbf{x}) := W(\mu_N[\mathbf{x}]), \quad \forall \mathbf{x} \in \mathcal{X}^N.$$

We recall from [16] that the value function  $V$  of the CMKV-MDP is  $\gamma$ -Hölder:

$$|V(\mu) - V(\mu')| \leq K_\star (\mathcal{W}_d(\mu, \mu'))^\gamma, \quad \forall \mu, \mu' \in \mathcal{P}(\mathcal{X}), \quad (3.2)$$

for some constant  $K_\star$  depending on  $K_F$ ,  $\beta$  and  $\Delta_{\mathcal{X}}$ .

**Lemma 3.1** *There exists some positive constant  $C$  such that for all  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)$ ,  $\mathbf{a} \in A^N$ ,  $\mathbf{x} \in \mathcal{X}^N$  and  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ ,*

$$|\widetilde{\mathbb{T}^{\mathbf{a}} V}(\mathbf{x}) - \mathbb{T}_N^{\mathbf{a}} \check{V}(\mathbf{x})| \leq C \left[ (\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]))^\gamma + M_N^\gamma \right].$$

**Proof.** For any  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)$ , and  $\mathbf{a} = (a^i)_{i \in [1, N]} \in A^N$ , we have

$$\begin{aligned} & \widetilde{\mathbb{T}^{\mathbf{a}} V}(\mathbf{x}) - \mathbb{T}_N^{\mathbf{a}} \check{V}(\mathbf{x}) \\ &= \mathbb{E} \left[ f(\xi_{\mathbf{x}}, a(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U))) - \frac{1}{N} \sum_{i=1}^N f(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}]) \right] \\ &+ \beta \mathbb{E} \left[ V(\mathbb{P}_{F(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1, \varepsilon_1^0)}^0) - V\left(\frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)}\right) \right]. \quad (3.3) \end{aligned}$$

We write

$$\mathbb{E} \left[ f(\xi_{\mathbf{x}}, a(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U))) \right] - \frac{1}{N} \sum_{i=1}^N f(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}]) = \hat{f}(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U))) - \hat{f}(\mu_N[\mathbf{x}, \mathbf{a}]),$$

where  $\hat{f}(\mu) = \int f(x', a', \mu) \mu(dx', da')$  for all  $\mu \in \mathcal{P}(\mathcal{X} \times A)$ . Notice that for  $\mu, \mu' \in \mathcal{P}(\mathcal{X} \times A)$ , we have

$$\begin{aligned} \hat{f}(\mu) - \hat{f}(\mu') &= \int f(x', a', \mu) (\mu - \mu')(dx', da') + \int (f(x', a', \mu) - f(x', a', \mu')) \mu'(dx', da') \\ &\leq K_f \mathcal{W}_d(\mu, \mu') + K_f \mathcal{W}_d(\mu, \mu') = 2K_f \mathcal{W}_d(\mu, \mu'), \end{aligned}$$

from the Kantorovich-Rubinstein dual representation (2.1) and  $(\mathbf{Hf}_{\text{ip}})$ . It follows that

$$\begin{aligned} &\left| \mathbb{E} \left[ f(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U))) \right] - \frac{1}{N} \sum_{i=1}^N f(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}]) \right| \\ &\leq 2K_f \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]). \end{aligned} \quad (3.4)$$

Let us next focus on the second term in (3.3). As  $V$  is  $\gamma$ -Hölder with constant factor  $K_*$ , we have

$$\begin{aligned} &\left| \mathbb{E} \left[ V \left( \mathbb{P}_{F(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1, \varepsilon_1^0)}^0 \right) - V \left( \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \right| \\ &\leq K_* \mathbb{E} \left[ \mathcal{W}_d \left( \mathbb{P}_{F(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^i, \varepsilon_1^0)}^0, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right]^\gamma, \end{aligned} \quad (3.5)$$

by Jensen's inequality. Let  $(\xi^i, U_0^i)_{i \in [1, N]}$  be  $N$  i.i.d. random variables, independent of  $\varepsilon_1$ , such that  $(\xi^i, U_0^i) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ ,  $i \in [1, N]$ . For any i.i.d. random variables  $(\tilde{\varepsilon}_1^i)_{i \in [1, N]}$  such that

$$((\xi^i, U_0^i, \tilde{\varepsilon}_1^i)_{i \in [1, N]}, \varepsilon_1^0) \stackrel{d}{=} ((\xi^i, U_0^i, \varepsilon_1^i)_{i \in [1, N]}, \varepsilon_1^0), \quad (3.6)$$

we have

$$\begin{aligned} &\mathbb{E} \left[ \mathcal{W}_d \left( \mathbb{P}_{F(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^i, \varepsilon_1^0)}^0, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \\ &\leq \mathbb{E} \left[ \mathcal{W}_d \left( \mathbb{P}_{F(\xi^i, \mathbf{a}(\xi^i, U_0^i), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^i, \varepsilon_1^0)}^0, \frac{1}{N} \sum_{i=1}^N \delta_{F(\xi^i, \mathbf{a}(\xi^i, U_0^i), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \tilde{\varepsilon}_1^i, \varepsilon_1^0)} \right) \right] \\ &\quad + \mathbb{E} \left[ \mathcal{W}_d \left( \frac{1}{N} \sum_{i=1}^N \delta_{F(\xi^i, \mathbf{a}(\xi^i, U_0^i), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \tilde{\varepsilon}_1^i, \varepsilon_1^0)}, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \\ &\leq M_N + \mathbb{E} \left[ \mathcal{W}_d \left( \frac{1}{N} \sum_{i=1}^N \delta_{F(\xi^i, \mathbf{a}(\xi^i, U_0^i), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \tilde{\varepsilon}_1^i, \varepsilon_1^0)}, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right], \end{aligned} \quad (3.7)$$

by definition of  $M_N$  in (2.3). Let us now consider the random permutation  $\sigma^{(\xi^i, \mathbf{a}(\xi^i, U_0^i))_{i \in [1, N]}, (x^i, a^i)_{i \in [1, N]}}$  defined in Definition 3.1 that we shall, to simplify notations, simply denote by  $\sigma$ . Notice that as  $(\xi^i, \mathbf{a}(\xi^i, U_0^i))_{i \in [1, N]} \perp (\varepsilon_1^i)_{i \in [1, N]}$ , we clearly see that  $(\tilde{\varepsilon}_1^i)_{i \in [1, N]} := (\varepsilon_1^{(\sigma^{-1})^i})_{i \in [1, N]}$  satisfies the required condition (3.6). Therefore the above relation applies to  $(\tilde{\varepsilon}_1^i)_{i \in [1, N]} = (\varepsilon_1^{(\sigma^{-1})^i})_{i \in [1, N]}$ .

For such  $(\tilde{\varepsilon}_1^i)_{i \in [1, N]}$ , we get

$$\begin{aligned}
& \mathbb{E} \left[ \mathcal{W}_d \left( \frac{1}{N} \sum_{i=1}^N \delta_{F(\xi^i, \mathbf{a}(\xi^i, U_0^i), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^{(\sigma^{-1})^i}, \varepsilon_1^0)}, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, \mathbf{a}^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \\
&= \mathbb{E} \left[ \mathcal{W}_d \left( \frac{1}{N} \sum_{i=1}^N \delta_{F(\xi^{\sigma^i}, \mathbf{a}(\xi^{\sigma^i}, U_0^{\sigma^i}), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^i, \varepsilon_1^0)}, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, \mathbf{a}^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \\
&\leq \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ d(F(\xi^{\sigma^i}, \mathbf{a}(\xi^{\sigma^i}, U_0^{\sigma^i}), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^i, \varepsilon_1^0), F(x^i, \mathbf{a}^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)) \right] \\
&\leq K_F \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \mathbf{d}((\xi^{\sigma^i}, \mathbf{a}(\xi^{\sigma^i}, U_0^{\sigma^i})), (x^i, \mathbf{a}^i)) + \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right] \\
&= K_F \mathbb{E} \left[ \mathcal{W}_d \left( \frac{1}{N} \sum_{i=1}^N \delta_{(\xi^i, \mathbf{a}(\xi^i, U_0^i))}, \mu_N[\mathbf{x}, \mathbf{a}] \right) + \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right] \\
&\leq K_F \left( M_N + 2 \mathbb{E} \left[ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right] \right),
\end{aligned}$$

where the first inequality comes from (2.2), the second one is derived by conditioning w.r.t.  $((\xi^i, U_0^i)_{i \in [1, N]}, \varepsilon_1^0)$  and using the regularity in expectation of  $F$  in  $(\mathbf{HF}_{\text{lip}})$ , the last equality holds true by definition of the permutation  $\sigma$  realizing the optimal coupling (3.1), and the last inequality from the definition of  $M_N$ . Recalling (3.7), we then have

$$\begin{aligned}
& \mathbb{E} \left[ \mathcal{W}_d \left( \mathbb{P}_F^0(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1^i, \varepsilon_1^0)}, \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, \mathbf{a}^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \\
&\leq (1 + K_F) M_N + 2K_F \mathbb{E} \left[ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right]
\end{aligned}$$

which implies by (3.5)

$$\begin{aligned}
& \mathbb{E} \left[ V \left( \mathbb{P}_F^0(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U), \mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \varepsilon_1, \varepsilon_1^0) \right) - V \left( \frac{1}{N} \sum_{i=1}^N \delta_{F(x^i, \mathbf{a}^i, \mu_N[\mathbf{x}, \mathbf{a}], \varepsilon_1^i, \varepsilon_1^0)} \right) \right] \\
&\leq K_\star \left( (1 + K_F) M_N + 2K_F \mathbb{E} \left[ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right] \right)^\gamma.
\end{aligned}$$

Together with (3.4), and plugging into (3.3), we obtain finally

$$\begin{aligned}
& \left| \widetilde{\mathbb{T}}^a \overline{V}(\mathbf{x}) - \mathbb{T}_N^{\mathbf{a}} \check{V}(\mathbf{x}) \right| \\
&\leq 2K_f \mathbb{E} \left[ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right] + K_\star \left( (1 + K_F) M_N + 2K_F \mathbb{E} \left[ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right] \right)^\gamma \\
&\leq C \left\{ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) + \left( \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \right)^\gamma + M_N^\gamma \right\}
\end{aligned}$$

(recall that  $\gamma \leq 1$ ), for some constant  $C$  depending only on  $K_\star$ ,  $K_f$ ,  $K_F$ , where we also use the fact that  $\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}])$  is bounded by a constant depending on the diameter of the compact set  $\mathcal{X} \times A$ . This ends the proof.  $\square$

### 3.2 Proof of Theorem 2.1

Lemma 3.1 means that given  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)$ ,  $\mathbf{a} \in A^N$ , and for  $\mathbf{x} \in \mathcal{X}^N$ , the Wasserstein distance between the distribution law of  $(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U))$  (where  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ ), and

the empirical measure  $\mu_N[\mathbf{x}, \mathbf{a}]$  is small (and  $N$  large), then  $\mathbb{T}^a V \simeq \mathbb{T}_N^{\mathbf{a}} \check{V}$ . It is thus natural to look for suitable choices of  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)$ ,  $\mathbf{a} \in A^N$  so that the above Wasserstein distance is as small as possible. This is quantified in the following result.

**Lemma 3.2** *Fix  $\mathbf{x} \in \mathcal{X}^N$ . Then, for any  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)$ , there exists  $\mathbf{a}^{\mathbf{a}} \in A^N$  such that*

$$\mathcal{W}_{\mathbf{d}}(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}^{\mathbf{a}}]) \leq 2M_N,$$

where  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ . Conversely, for any  $\mathbf{a} \in A^N$ , there exists  $\mathbf{a}^{\mathbf{a}} \in L^0(\mathcal{X} \times [0, 1]; A)$  such that

$$\mathcal{W}_{\mathbf{d}}(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}^{\mathbf{a}}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) = 0.$$

**Proof.** Fix  $\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)$ . Let us consider  $\xi = (\xi^i)_{i \in [1, N]}$  i.i.d. with common distribution  $\mu_N[\mathbf{x}]$ , independent from  $U_0 = (U_0^i)_{i \in [1, N]}$  i.i.d.  $\sim \mathcal{U}([0, 1])$ . We have

$$\begin{aligned} & \mathbb{E}\left[\mathcal{W}_{\mathbf{d}}\left(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \frac{1}{N} \sum_{i=1}^N \delta_{x^i, \mathbf{a}(\xi^{\sigma_i^{\xi, \mathbf{x}}}, U_0^i)}\right)\right] \\ & \leq \mathbb{E}\left[\mathcal{W}\left(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \frac{1}{N} \sum_{i=1}^N \delta_{\xi^{\sigma_i^{\xi, \mathbf{x}}}, \mathbf{a}(\xi^{\sigma_i^{\xi, \mathbf{x}}}, U_0^i)}\right) + \mathcal{W}_{\mathbf{d}}\left(\frac{1}{N} \sum_{i=1}^N \delta_{\xi^{\sigma_i^{\xi, \mathbf{x}}}, \mathbf{a}(\xi^{\sigma_i^{\xi, \mathbf{x}}}, U_0^i)}, \frac{1}{N} \sum_{i=1}^N \delta_{x^i, \mathbf{a}(\xi^{\sigma_i^{\xi, \mathbf{x}}}, U_0^i)}\right)\right] \\ & \leq M_N + \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N d(\xi^{\sigma_i^{\xi, \mathbf{x}}}, x^i)\right] \leq 2M_N, \end{aligned}$$

where we used the definition of  $M_N$  and (2.2) in the second inequality, and definition of  $\sigma^{\xi, \mathbf{x}}$  in the last inequality. It follows that

$$\mathbb{P}\left[\mathcal{W}_{\mathbf{d}}\left(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \frac{1}{N} \sum_{i=1}^N \delta_{x^i, \mathbf{a}(\xi^{\sigma_i^{\xi, \mathbf{x}}}, U_0^i)}\right) \leq 2M_N\right] > 0,$$

which implies that there exists a vector  $\mathbf{a} \in A^N$  such that

$$\mathcal{W}_{\mathbf{d}}(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) \leq 2M_N.$$

On the other hand, given such an  $\mathbf{a} \in A^N$ , there clearly exists  $\mathbf{a}^{\mathbf{a}} \in L^0(\mathcal{X} \times [0, 1]; A)$  such that  $\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}^{\mathbf{a}}(\xi_{\mathbf{x}}, U)) = \mu_N[\mathbf{x}, \mathbf{a}]$ : indeed, by considering  $(\tilde{\xi}, \tilde{\alpha}) \sim \mu_N[\mathbf{x}, \mathbf{a}]$ , it suffices to choose  $\mathbf{a}^{\mathbf{a}}$  as a kernel for simulating the conditional distribution of  $\tilde{\alpha}$  knowing  $\tilde{\xi}$ . We then have

$$\mathcal{W}(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}^{\mathbf{a}}(\xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]) = 0. \quad \square$$

By combining the general comparison of Bellman operators in Lemma 3.1 with the coupling result in Lemma 3.2, we can now prove the propagation of chaos of value functions.

**Proof of Theorem 2.1.** From the fixed point equation for  $V$  with Bellman operator  $\mathcal{T}$  in (2.4), we have

$$\begin{aligned} \check{V}(\mathbf{x}) &= \widetilde{\mathcal{T}}V(\mathbf{x}) \\ &= \sup_{\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)} \widetilde{\mathbb{T}}^{\mathbf{a}}V(\mathbf{x}) \leq \sup_{\mathbf{a} \in L^0(\mathcal{X} \times [0, 1]; A)} \mathbb{T}_N^{\mathbf{a}}\check{V}(\mathbf{x}) + CM_N^\gamma \\ &\leq \mathcal{T}_N\check{V}(\mathbf{x}) + CM_N^\gamma, \end{aligned}$$

where we used Lemma 3.1 and Lemma 3.2 in the first inequality, and the definition of  $\mathcal{T}_N$  in the last one. Since  $V_N$  is a fixed point of  $\mathcal{T}_N$  (see Proposition B.2), we then have:

$$(\check{V} - V_N)(\mathbf{x}) \leq (\mathcal{T}_N \check{V} - \mathcal{T}_N V_N)(\mathbf{x}) + CM_N^\gamma,$$

and thus by definition of  $\mathcal{T}_N$ ,

$$(\check{V} - V_N)(\mathbf{x}) \leq \beta \sup_{\mathbf{x}' \in \mathcal{X}^N} (\check{V} - V_N)(\mathbf{x}') + CM_N^\gamma,$$

which implies

$$\sup_{\mathbf{x} \in \mathcal{X}^N} (\check{V} - V_N)(\mathbf{x}) \leq CM_N^\gamma.$$

Likewise, by Lemma 3.1 and Lemma 3.2, we have

$$\begin{aligned} \check{V}(\mathbf{x}) &= \widetilde{\mathcal{T}V}(\mathbf{x}) = \sup_{\mathbf{a} \in L^0(\mathcal{X} \times [0,1]; A)} \widetilde{\mathbb{T}^{\mathbf{a}}V}(\mathbf{x}) \geq \sup_{\mathbf{a} \in A^N} \widetilde{\mathbb{T}^{\mathbf{a}^{\mathbf{a}}}V}(\mathbf{x}) \\ &\geq \sup_{\mathbf{a} \in A^N} \mathbb{T}_N^{\mathbf{a}} \check{V}(\mathbf{x}) - CM_N^\gamma = \mathcal{T}_N \check{V}(\mathbf{x}) - CM_N^\gamma, \end{aligned}$$

and using the fact that  $V_N$  is a fixed point of  $\mathcal{T}_N$ , we obtain similarly

$$\sup_{\mathbf{x} \in \mathcal{X}^N} (V_N - \check{V})(\mathbf{x}) \leq CM_N^\gamma,$$

which concludes the proof.

### 3.3 Proof of Theorem 2.2

We start with a general result estimating the efficiency of a feedback policy for the  $N$ -agent MDP by “comparing” it to an  $\epsilon$ -optimal randomized feedback policy for the CMKV- MDP.

**Lemma 3.3** *Let  $\mathbf{a}_\epsilon$  be an  $\epsilon$ -optimal randomized feedback policy for the CMKV-MDP, and  $\mathbf{a} \in A^N$ . Then, there exists some positive constant  $C$  (depending only on  $\Delta_{\mathcal{X} \times A}$ ,  $\beta$ ,  $K_F$ ,  $K_f$ ) such that for all  $\mathbf{x} \in \mathcal{X}^N$ ,*

$$\mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x}) \geq V_N(\mathbf{x}) - \epsilon - C[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}])^\gamma + M_N^\gamma].$$

**Proof.** Fix  $\mathbf{x} \in \mathcal{X}^N$ ,  $\mathbf{a} \in A^N$ , and define  $\mathbf{a}_\epsilon \in L^0(\mathcal{X} \times [0,1]; A)$  by  $\mathbf{a}_\epsilon(x, u) = \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], x, u)$  for  $x \in \mathcal{X}$ ,  $u \in [0,1]$ . By Theorem 2.1 and the  $\beta$ -contracting property of  $\mathbb{T}_N^{\mathbf{a}}$ , we have

$$|\mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x}) - \mathbb{T}_N^{\mathbf{a}} \check{V}(\mathbf{x})| \leq \beta \|V_N(\mathbf{x}) - \check{V}(\mathbf{x})\|_{\mathcal{X}^N} \leq CM_N^\gamma,$$

and so

$$\mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x}) \geq \mathbb{T}_N^{\mathbf{a}} \check{V}(\mathbf{x}) - CM_N^\gamma.$$

Together with Lemma 3.1, this yields

$$\mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x}) \geq \widetilde{\mathbb{T}^{\mathbf{a}_\epsilon}V}(\mathbf{x}) - C[\mathcal{W}_d(\mathcal{L}(\xi, \mathbf{a}_\epsilon(\xi, U)), \mu_N[\mathbf{x}, \mathbf{a}])^\gamma + M_N^\gamma]. \quad (3.8)$$

Denote by  $\alpha^\epsilon$  the randomized feedback control associated via (2.5) to the randomized feedback policy  $\mathbf{a}_\epsilon$ . Then, notice that the gain functional  $V^{\alpha^\epsilon}(\xi)$  depends on  $\xi$  only through its law  $\mu = \mathcal{L}(\xi)$ , and we set  $V^{\alpha^\epsilon}(\mu) = V^{\alpha^\epsilon}(\xi)$  when  $\xi \sim \mu$ . Since  $V \geq V^{\alpha^\epsilon}$ , and by the monotonicity of  $\mathbb{T}^{\mathbf{a}_\epsilon}$ , we have

$$\mathbb{T}^{\mathbf{a}_\epsilon} V(\mu_N[\mathbf{x}]) \geq \mathbb{T}^{\mathbf{a}_\epsilon} V^{\alpha^\epsilon}(\mu_N[\mathbf{x}]) = V^{\alpha^\epsilon}(\mu_N[\mathbf{x}]) \geq V(\mu_N[\mathbf{x}]) - C\epsilon$$

by recalling that  $V^{\alpha^\epsilon}$  is a fixed point of  $\mathbb{T}^{\mathbf{a}_\epsilon}$ , and using the fact that  $\mathbf{a}_\epsilon$  is an  $\epsilon$ -optimal randomized feedback policy for the CMKV-MDP. From Theorem 2.1, this implies that

$$\widetilde{\mathbb{T}^{\mathbf{a}_\epsilon} V}(\mathbf{x}) \geq V_N(\mathbf{x}) - C[\epsilon + M_N^\gamma],$$

which proved the required result when combined with (3.8).  $\square$

Let us denote by  $L^0(\mathcal{X}^N; A^N)$  the set of measurable functions from  $\mathcal{X}^N$  into  $A^N$ . Given a feedback policy  $\pi \in L^0(\mathcal{X}^N; A^N)$  for the  $N$ -agent problem, the associated feedback control is the unique control  $\alpha^\pi$  defined by  $\alpha_t^\pi = \pi(\mathbf{X}_t)$ ,  $t \in \mathbb{N}$ . By misuse of notation, we denote  $V_N^\pi = V_N^{\alpha^\pi}$ . Let us then introduce the operator  $\mathcal{T}_N^\pi$  on  $L_m^\infty(\mathcal{X}^N)$ , defined by

$$\mathcal{T}_N^\pi W(\mathbf{x}) = \mathbf{f}(\mathbf{x}, \pi(\mathbf{x})) + \beta \mathbb{E}[W(\mathbf{F}(\mathbf{x}, \pi(\mathbf{x}), \varepsilon_1))], \quad \mathbf{x} \in \mathcal{X}^N.$$

**Proposition 3.1** *Let  $\mathbf{a}_\epsilon$  be an  $\epsilon$ -optimal randomized feedback policy for the CMKV-MDP, and consider any feedback policy  $\pi$  for the  $N$ -agent MDP. Then, the feedback control  $\alpha^\pi$  is*

$$\mathcal{O}(\epsilon + \sup_{\mathbf{x} \in \mathcal{X}^N} \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi(\mathbf{x})])^\gamma + M_N^\gamma)\text{-optimal for } V_N(\mathbf{x}_0),$$

where  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ , namely

$$V_N(\mathbf{x}_0) - C[\epsilon + \sup_{\mathbf{x} \in \mathcal{X}^N} \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi(\mathbf{x})])^\gamma + M_N^\gamma] \leq V_N^\pi(\mathbf{x}_0).$$

**Proof.** Fix  $\mathbf{x} \in \mathcal{X}^N$ , and let  $\mathbf{a} = \pi(\mathbf{x}) \in A^N$ . By definition, we have  $\mathcal{T}_N^\pi V_N(\mathbf{x}) = \mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x})$ . By Lemma 3.3, we thus have

$$\mathcal{T}_N^\pi V_N(\mathbf{x}) \geq V_N(\mathbf{x}) - \epsilon - C[\sup_{\mathbf{x} \in \mathcal{X}^N} \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}])^\gamma + M_N^\gamma],$$

and we conclude by the verification result in Lemma B.5.  $\square$

Proposition 3.1 has an important implication: it means that a feedback policy  $\pi$  for the  $N$ -agent MDP yields the better performance whenever it assigns for each state  $\mathbf{x}$  the action  $\pi(\mathbf{x})$  that achieves the minimum of

$$\mathbf{a} \in A^N \mapsto \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]).$$

Let us check that one can choose a measurable version of this argmin.

**Lemma 3.4** *Let  $\mathbf{a} \in L^0(\mathcal{P}(\mathcal{X}) \times \mathcal{X} \times [0, 1]; A)$ . Then, there exists a measurable function  $\pi^* : \mathcal{X}^N \rightarrow A^N$  such that*

$$\pi^*(\mathbf{x}) \in \operatorname{argmin}_{\mathbf{a} \in A^N} \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]), \quad \mathbf{x} \in \mathcal{X}^N.$$



**Proof.** Notice that the function

$$h(\mathbf{x}, \mathbf{a}) := \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}])$$

is such that for  $\mathbf{a} \in A^N$ ,  $h(\cdot, \mathbf{a})$  is measurable, and for  $\mathbf{x} \in \mathcal{X}^N$ ,  $h(\mathbf{x}, \cdot)$  is continuous. Let us then show that one can measurably select  $\operatorname{argmin} h(\mathbf{x}, \mathbf{a})$  w.r.t.  $\mathbf{x}$ . Consider a dense sequence  $(\mathbf{a}_n)_{n \in \mathbb{N}} \subset A^N$  (its existence is guaranteed by the fact that  $A^N$  is a compact metric space), and define by recursion the sequence of measurable functions  $\pi_n : \mathcal{X}^N \rightarrow A^N$  as

$$\begin{aligned} \pi_0(\mathbf{x}) &= \mathbf{a}_0 \\ \pi_{n+1}(\mathbf{x}) &= \begin{cases} \pi_n(\mathbf{x}) & \text{if } h(\mathbf{x}, \pi_n(\mathbf{x})) \leq h(\mathbf{x}, \mathbf{a}_{n+1}) \\ \mathbf{a}_{n+1} & \text{else.} \end{cases} \end{aligned}$$

The measurability of  $\pi_n$  is easily established by induction on  $n$ : For  $n = 0$ , it is clear. Assuming that  $\pi_n$  is measurable, and denoting

$$g_n(\mathbf{x}) = h(\mathbf{x}, \pi_n(\mathbf{x})) - h(\mathbf{x}, \mathbf{a}_{n+1}), \forall \mathbf{x} \in \mathcal{X}^N,$$

notice that for any measurable set  $B \subset A^N$ , we have

$$[\pi_{n+1}]^{-1}(B) = \begin{cases} [\pi_n]^{-1}(B) \cap g_n^{-1}(\mathbb{R}_-) & \text{if } \mathbf{a}_{n+1} \notin B, \\ ([\pi_n]^{-1}(B) \cap g_n^{-1}(\mathbb{R}_-)) \cup g_n^{-1}(\mathbb{R}_+) & \text{if } \mathbf{a}_{n+1} \in B, \end{cases}$$

which is clearly a measurable set, and proves the induction. Then, let us consider an embedding  $\phi : A^N \rightarrow [0, 1]$  such that  $\phi$  and  $\phi^{-1}$  are uniformly continuous (see Lemma C.2 in [16]). Then,  $(\phi \circ \pi_n)_{n \in \mathbb{N}}$  denotes a sequence of measurable functions from  $\mathcal{X}^N$  to  $\phi(A^N) \subset [0, 1]$ . It is well known that the function  $\liminf_{n \in \mathbb{N}} (\phi \circ \pi_n)$  is then measurable from  $\mathcal{X}^N$  to  $\phi(A^N)$  (we here use the fact that  $\phi$  is continuous and  $\phi(A^N)$  is closed, which ensures that the  $\liminf$  takes its values in  $\phi(A^N)$ ). Finally, let us denote  $\pi^* : \mathcal{X}^N \rightarrow A^N$  defined by

$$\pi^* = \phi^{-1} \circ \liminf_{n \in \mathbb{N}} (\phi \circ \pi_n).$$

$\pi^*$  is then measurable by composition. Furthermore, for any  $\mathbf{x} \in \mathcal{X}^N$ ,  $\phi \circ \pi^*(\mathbf{x}) = \liminf_{n \in \mathbb{N}} (\phi \circ \pi_n(\mathbf{x}))$  is an accumulation point of the sequence  $(\phi \circ \pi_n(\mathbf{x}))_{n \in \mathbb{N}}$ , which implies, by continuity of  $\phi^{-1}$ , that  $\pi^*(\mathbf{x})$  is an accumulation point of  $(\pi_n(\mathbf{x}))_{n \in \mathbb{N}}$ . Given the definition of  $\pi_n(\mathbf{x})$ , it is clear by induction that for any  $n \in \mathbb{N}$ ,  $h(\mathbf{x}, \pi_n(\mathbf{x})) \leq \min_{m \leq n} h(\mathbf{x}, \pi_m^*)$ , and thus  $h(\mathbf{x}, \pi^*(\mathbf{x})) \leq \min_{n \in \mathbb{N}} h(\mathbf{x}, \pi_n)$ . By density of  $(\mathbf{a}^n)_{n \in \mathbb{N}}$  and by continuity of  $h(\mathbf{x}, \cdot)$ , this implies that  $h(\mathbf{x}, \pi^*(\mathbf{x})) = \min_{\mathbf{a} \in A^N} h(\mathbf{x}, \mathbf{a})$  for all  $\mathbf{x} \in \mathcal{X}^N$ , i.e.  $\pi^*(\mathbf{x}) \in \operatorname{argmin}_{\mathbf{a} \in A^N} h(\mathbf{x}, \mathbf{a})$ . We conclude that  $\pi^*$  is thus a measurable selection of  $\operatorname{argmin} h(\cdot, \mathbf{a})$ .  $\square$

By Lemma 3.4, there exists a randomized feedback policy  $\pi^{\mathbf{a}_\epsilon, N}$  s.t.

$$\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi^{\mathbf{a}_\epsilon, N}]) = \inf_{\mathbf{a} \in A^N} \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\epsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \mathbf{a}]),$$

and the r.h.s. of the above equality is bounded by  $2M_N$  from Lemma 3.2. Together with Proposition 3.1, this proves Theorem 2.2.

### 3.4 Proof of Theorem 2.3

Given a randomized feedback policy  $\pi_r \in L^0(\mathcal{X}^N \times [0, 1]^N; A^N)$ , the set of measurable functions from  $\mathcal{X}^N \times [0, 1]^N$  into  $A^N$ , the associated feedback control is the unique control  $\alpha^\pi$  given by  $\alpha_t^\pi = \pi_r(\mathbf{X}_t, \mathbf{U}_t)$ ,  $t \in \mathbb{N}$ , where  $\{\mathbf{U}_t = (U_t^i)_{i \in [1, N]}\}$ ,  $t \in \mathbb{N}$  is a family of mutually i.i.d. uniform random variables on  $[0, 1]$ , independent of  $\mathcal{G}$ ,  $\varepsilon$ . By misuse of notation, we denote  $V_N^\pi = V_N^{\alpha^\pi}$ . For  $\pi_r \in L^0(\mathcal{X}^N \times [0, 1]^N; A^N)$ , we introduce the operator  $\mathcal{T}_N^{\pi_r}$  on  $L_m^\infty(\mathcal{X}^N)$ , defined by

$$\mathcal{T}_N^{\pi_r} W(\mathbf{x}) := \mathbb{E}[\mathbf{f}(\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)) + \beta W(\mathbf{F}(\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0), \varepsilon_1))], \quad \forall \mathbf{x} \in \mathcal{X}^N,$$

where  $\mathbf{U}_0 = (U_0^i)_{i \in [1, N]}$  is a family of i.i.d.  $\sim \mathcal{U}([0, 1])$ , independent of  $\mathcal{G}$ ,  $\varepsilon$ .

We adapt Proposition 3.1 to the case of randomized feedback policies.

**Proposition 3.2** *Let  $\mathbf{a}_\varepsilon$  be an  $\varepsilon$ -optimal randomized feedback policy for the CMKV-MDP, and consider any feedback policy  $\pi_r$  for the  $N$ -agent MDP. Then, the feedback control  $\alpha^{\pi_r}$  is*

$$\mathcal{O}(\varepsilon + \sup_{\mathbf{x} \in \mathcal{X}^N} \mathbb{E}[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\varepsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)])]^\gamma + M_N^\gamma)\text{-optimal for } V_N(\mathbf{x}_0),$$

namely

$$V_N(\mathbf{x}_0) - C \left( \varepsilon + \sup_{\mathbf{x} \in \mathcal{X}^N} \mathbb{E}[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\varepsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)])]^\gamma + M_N^\gamma \right) \leq V_N^\pi(\mathbf{x}_0).$$

Here  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ , and  $\mathbf{U}_0 = (U_0^i)_{i \in [1, N]}$  is a family of i.i.d.  $\sim \mathcal{U}([0, 1])$ , independent of  $\varepsilon$ .

**Proof.** Fix  $\mathbf{x} \in \mathcal{X}^N$ , and let  $\mathbf{a} = \pi_r(\mathbf{x}, \mathbf{U}_0)$  be the random variable valued in  $A^N$ . By definition, we have  $\mathcal{T}_N^{\pi_r} V_N(\mathbf{x}) = \mathbb{E}[\mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x})]$ . By Lemma 3.3, we have

$$\mathbb{T}_N^{\mathbf{a}} V_N(\mathbf{x}) \geq V_N(\mathbf{x}) - \varepsilon - C[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\varepsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)])]^\gamma + M_N^\gamma].$$

Taking the expectation, and by Jensen's inequality, we then get

$$\mathcal{T}_N^{\pi_r} V_N(\mathbf{x}) \geq V_N(\mathbf{x}) - \varepsilon - C \left( \sup_{\mathbf{x} \in \mathcal{X}^N} \mathbb{E}[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\varepsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)])]^\gamma + M_N^\gamma \right),$$

and we conclude by the verification result in Lemma B.5.  $\square$

Compared to Proposition 3.1, Proposition 3.2 means that with a randomized feedback policy  $\pi_r$ , one can obtain a ‘‘good’’ performance whenever it produces empirical state-action distributions that are close the theoretical state-action distribution generated by  $\mathbf{a}_\varepsilon$  on average, i.e., that makes the quantity

$$\mathbb{E}[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\varepsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)])]$$

as small as possible. More precisely, if we can design a randomized policy  $\pi_r$  such that

$$\mathbb{E}[\mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_\varepsilon(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r(\mathbf{x}, \mathbf{U}_0)])] \leq CM_N,$$

then by Proposition 3.2, this will prove the statement of Theorem 2.3. The next result shows how it can be achieved.

**Lemma 3.5** Let  $\mathbf{a} : \mathcal{P}(\mathcal{X}) \times \mathcal{X} \times [0, 1] \rightarrow A$  be any (if it exists) randomized feedback policy for the CMKV-MDP such that

$$\mathbb{E}[d_A(\mathbf{a}(\mu, x, U), \mathbf{a}(\mu, x', U))] \leq Kd(x, x'), \quad \forall \mu \in \mathcal{P}(\mathcal{X}), x, x' \in \mathcal{X}, \quad (3.9)$$

(here  $U \sim \mathcal{U}[0, 1]$ ) for some positive constant  $K$ . Consider the randomized feedback policy for the  $N$ -agent MDP defined by

$$\pi_r^{\mathbf{a}, N}(\mathbf{x}, \mathbf{u}) = \left( \mathbf{a}(\mu_N[\mathbf{x}], x^i, u^i) \right)_{i \in [1, N]}, \quad \mathbf{x} = (x^i)_{i \in [1, N]} \in \mathcal{X}^N, \quad \mathbf{u} = (u^i)_{i \in [1, N]} \in [0, 1]^N.$$

Then,

$$\mathbb{E} \left[ \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r^{\mathbf{a}, N}(\mathbf{x}, \mathbf{U}_0)]) \right] \leq (2 + K)M_N,$$

where  $(\xi_{\mathbf{x}}, U) \sim \mu_N[\mathbf{x}] \otimes \mathcal{U}([0, 1])$ .

**Proof.** Fix  $\mathbf{x} \in \mathcal{X}^N$ , and set  $\mathbf{a}_{\mathbf{x}}(x, u) = \mathbf{a}(\mu_N[\mathbf{x}], x, u)$  for  $(x, u) \in \mathcal{X} \times [0, 1]$ . Let us consider a family  $\boldsymbol{\xi} = (\xi^i)_{i \in [1, N]}$  of  $N$  i.i.d. random variables such that  $\xi^i \sim \mu_N[\mathbf{x}]$ , and independent of  $\mathbf{U}_0$ . Let us consider  $\sigma^{\boldsymbol{\xi}, \mathbf{x}}$ , the optimal permutation defined in Definition 3.1 between  $\boldsymbol{\xi}$  and  $\mathbf{x}$ . We have

$$\begin{aligned} & \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r^{\mathbf{a}, N}(\mathbf{x}, \mathbf{U}_0)]) \\ &= \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_{\mathbf{x}}(\xi_{\mathbf{x}}, U)), \frac{1}{N} \sum_{i=1}^N \delta_{x^i, \mathbf{a}_{\mathbf{x}}(x^i, U_0^i)}) \\ &\leq \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_{\mathbf{x}}(\xi_{\mathbf{x}}, U)), \frac{1}{N} \sum_{i=1}^N \delta_{\xi^{\sigma_i^{\boldsymbol{\xi}, \mathbf{x}}}, \mathbf{a}_{\mathbf{x}}(\xi^{\sigma_i^{\boldsymbol{\xi}, \mathbf{x}}}, U_0^i)}) + \mathcal{W}_d(\frac{1}{N} \sum_{i=1}^N \delta_{\xi^{\sigma_i^{\boldsymbol{\xi}, \mathbf{x}}}, \mathbf{a}_{\mathbf{x}}(\xi^{\sigma_i^{\boldsymbol{\xi}, \mathbf{x}}}, U_0^i)}, \frac{1}{N} \sum_{i=1}^N \delta_{x^i, \mathbf{a}_{\mathbf{x}}(x^i, U_0^i)}) \\ &\leq \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}_{\mathbf{x}}(\xi_{\mathbf{x}}, U)), \frac{1}{N} \sum_{i=1}^N \delta_{\xi^i, \mathbf{a}_{\mathbf{x}}(\xi^i, U_0^{(\sigma^{\boldsymbol{\xi}, \mathbf{x}})^{-1}_i})}) + \mathbf{d}_N((\boldsymbol{\xi}^{\sigma^{\boldsymbol{\xi}, \mathbf{x}}}, \pi_r^{\mathbf{a}, N}(\boldsymbol{\xi}^{\sigma^{\boldsymbol{\xi}, \mathbf{x}}}, \mathbf{U}_0)), (\mathbf{x}, \pi_r^{\mathbf{a}, N}(\mathbf{x}, \mathbf{U}_0))), \end{aligned}$$

where we set  $\boldsymbol{\xi}^{\sigma^{\boldsymbol{\xi}, \mathbf{x}}} = (\xi^{\sigma_i^{\boldsymbol{\xi}, \mathbf{x}}})_{i \in [1, N]}$ , and use (2.2) in the last inequality. Taking the expectation, we then obtain under condition (3.9)

$$\begin{aligned} & \mathcal{W}_d(\mathcal{L}(\xi_{\mathbf{x}}, \mathbf{a}(\mu_N[\mathbf{x}], \xi_{\mathbf{x}}, U)), \mu_N[\mathbf{x}, \pi_r^{\mathbf{a}, N}(\mathbf{x}, \mathbf{U}_0)]) \\ &\leq M_N + (1 + K)\mathbb{E}[\mathbf{d}_N(\boldsymbol{\xi}^{\sigma^{\boldsymbol{\xi}, \mathbf{x}}}, \mathbf{x})] \\ &= M_N + (1 + K)\mathbb{E}[\mathcal{W}_d(\mu_N[\boldsymbol{\xi}], \mu_N[\mathbf{x}])] \leq (2 + K)M_N, \end{aligned}$$

where we use (3.1) in the last equality. This concludes the proof.  $\square$

We now apply Lemma 3.5 with an  $\varepsilon$ -optimal randomized feedback policy  $\mathbf{a}_\varepsilon$  for the CMKV-MDP, and combined with Proposition 3.2, this proves the required result in Theorem 2.3.

## A Existence of optimal randomized control for CMKV-MDP

Recall from Proposition 4.1 in [16] that the Bellman operator  $\mathcal{T}$  of the CMKV-MDP is written in the lifted form as

$$[\mathcal{T}W](\mu) = \sup_{\mathbf{a} \in \mathbf{A}} \left\{ \tilde{f}(\mu, \mathbf{a}) + \beta \mathbb{E}[W(\tilde{F}(\mu, \mathbf{a}, \varepsilon_1^0))] \right\}, \quad \mu \in \mathcal{P}(\mathcal{X}), \quad (\text{A.1})$$

for  $W \in L_m^\infty(\mathcal{P}(\mathcal{X}))$ , where  $\mathbf{A} = \mathcal{P}(\mathcal{X} \times A)$ ,  $\tilde{F}$  is the measurable function on  $\mathcal{P}(\mathcal{X}) \times \mathbf{A} \times E^0 \rightarrow \mathcal{P}(\mathcal{X})$  defined by

$$\tilde{F}(\mu, \mathbf{a}, e^0) = F(\cdot, \cdot, \mathbf{p}(\mu, \mathbf{a}), \cdot, e^0) \star (\mathbf{p}(\mu, \mathbf{a}) \otimes \mathcal{L}(\varepsilon_1)),$$

and  $\tilde{f}$  is the measurable function on  $\mathcal{P}(\mathcal{X}) \times \mathbf{A}$  defined by

$$\tilde{f}(\mu, \mathbf{a}) = \int_{\mathcal{X} \times A} f(x, a, \mathbf{p}(\mu, \mathbf{a})) \mathbf{p}(\mu, \mathbf{a})(dx, da).$$

Here  $\star$  is the pushforward measure notation,  $\mathbf{p}$  is a measurable coupling projection from  $\mathcal{P}(\mathcal{X}) \times \mathbf{A}$  into  $\mathbf{A}$ :  $\mathbf{p}(\mu, \mathbf{a}) = \mathbf{p}(\mu, \mathbf{p}(\mu, \mathbf{a}))$ , satisfying  $\text{pr}_1 \star \mathbf{p}(\mu, \mathbf{a}) = \mu$ , and  $\mathbf{p}(\mu, \mathbf{a}) = \mathbf{a}$  if  $\text{pr}_1 \star \mathbf{a} = \mu$  (where  $\text{pr}_1$  is the projection function on the first coordinate). Since  $\tilde{f}$  and  $\tilde{F}$  depend upon  $\mathbf{a}$  only through  $\mathbf{p}(\mu, \mathbf{a})$ , it is clear that the supremum in (A.1), for each  $\mu \in \mathcal{P}(\mathcal{X})$ , can be taken actually over the subset  $\Gamma_\mu := \{\mathbf{a} : (\mu, \mathbf{a}) \in \Gamma\} \subset \mathbf{A}$ , where  $\Gamma := \{(\mu, \mathbf{a}) \in \mathcal{P}(\mathcal{X}) \times \mathbf{A} : \text{pr}_1 \star \mathbf{a} = \mu\}$  is closed in  $\mathcal{P}(\mathcal{X}) \times \mathbf{A}$  from the continuity of  $\mathbf{a} \mapsto \text{pr}_1 \star \mathbf{a}$ . Moreover, since  $V$  is continuous (see (3.2)), it is straightforward to prove that

$$\begin{aligned} (\mu, \mathbf{a}) \in \Gamma &\mapsto \tilde{f}(\mu, \mathbf{a}) + \beta \mathbb{E} \left[ V(\tilde{F}(\mu, \mathbf{a}, \varepsilon_1^0)) \right] \\ &= \int_{\mathcal{X} \times A} f(x, a, \mathbf{a}) \mathbf{a}(dx, da) + \beta \mathbb{E} \left[ V(F(\cdot, \cdot, \mu, \cdot, e^0) \star (\mathbf{a} \otimes \mathcal{L}(\varepsilon_1))) \right] \end{aligned}$$

is continuous and thus upper continuous on  $\Gamma$ . Therefore, by [2], Proposition 7.33, there exists a measurable function  $\phi : \mathcal{P}(\mathcal{X}) \rightarrow \mathbf{A}$  whose graph is included in  $\Gamma$  and such that

$$\begin{aligned} \tilde{f}(\mu, \phi(\mu)) + \beta \mathbb{E} \left[ V(\tilde{F}(\mu, \phi(\mu), \varepsilon_1^0)) \right] &= \sup_{\mathbf{a} \in \Gamma_\mu} \left\{ \tilde{f}(\mu, \mathbf{a}) + \beta \mathbb{E} \left[ V(\tilde{F}(\mu, \mathbf{a}, \varepsilon_1^0)) \right] \right\}, \\ &= [TV](\mu) = V(\mu), \quad \forall \mu \in \mathcal{P}(\mathcal{X}), \end{aligned} \quad (\text{A.2})$$

where the last equality follows from the fixed point equation of  $V$ . By the universal disintegration theorem (see [13], Corollary 1.26), there exists  $\kappa : \mathcal{X} \times \mathcal{P}(\mathcal{X} \times A) \times \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(A)$  such that for all  $\mathbf{a} \in \mathcal{P}(\mathcal{X} \times A)$ ,  $\mu \in \mathcal{P}(\mathcal{X})$  with  $\text{pr}_1 \star \mathbf{a} = \mu$ , we have  $\mathbf{a} = \mu \hat{\otimes} \kappa(\cdot, \mathbf{a}, \mu)$  (where  $\hat{\otimes}$  denotes the probability-kernel product). Furthermore, by Blackwell-Dubins Lemma, there exists a measurable function  $\rho : \mathcal{P}(A) \times [0, 1] \rightarrow A$  such that for all  $\pi \in \mathcal{P}(A)$ , if  $U$  denotes a uniform random variable, then  $\rho(\pi, U) \sim \pi$ . We can then define the randomized feedback policy

$$\mathbf{a}_0(\mu, x, u) = \rho(\kappa(x, \phi(\mu), \mu), u),$$

which satisfies by construction  $\mathcal{L}(\xi, \mathbf{a}_0(\mu, \xi, U)) = \phi(\mu)$  for  $(\xi, U) \sim \mu \otimes \mathcal{U}([0, 1])$  so that

$$\begin{aligned} \tilde{f}(\mu, \phi(\mu)) &= \mathbb{E} \left[ f(\xi, \mathbf{a}_0(\mu, \xi, U), \mathcal{L}(\xi, \mathbf{a}_0(\mu, \xi, U))) \right] \\ \tilde{F}(\mu, \phi(\mu), \varepsilon_1^0) &= \mathbb{P}_{F(\xi, \mathbf{a}_0(\mu, \xi, U), \mathcal{L}(\xi, \mathbf{a}_0(\mu, \xi, U)), \varepsilon_1, \varepsilon_1^0)}^0. \end{aligned}$$

Recalling notation in (2.4), and by (A.2), this shows that

$$\mathbb{T}^{\mathbf{a}_0(\mu, \cdot)} V(\mu) = V(\mu).$$

According to the verification result (Proposition 4.3 in [16]), this ensures that that the randomized feedback control  $\alpha^0 \in \mathcal{A}$  defined by

$$\alpha_t^0 = \mathbf{a}_0(\mathbb{P}_{X_t}^0, X_t, U_t), \quad t \in \mathbb{N},$$

where  $(U_t)_{t \in \mathbb{N}}$  is an i.i.d. sequence of random variables,  $U_t \sim \mathcal{U}([0, 1])$ , independent of  $\xi_0 \sim \mu_0$ , and  $\varepsilon$ , is an optimal control for  $V(\mu_0)$ .

## B Bellman equation for the $N$ -agent MDP

In this section, we study and rigorously state properties on the Bellman equation for the  $N$ -agent problem, viewed as a MDP with state space  $\mathcal{X}^N$ , action space  $A^N$ , noise sequence  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_t)_{t \in \mathbb{N}^*}$  with  $\boldsymbol{\varepsilon}_t := ((\varepsilon_t^i)_{i \in \llbracket 1, N \rrbracket}, \varepsilon_t^0)$  valued in  $E^N \times E^0$ , state transition function

$$\mathbf{F}(\mathbf{x}, \mathbf{a}, \mathbf{e}) := \left( F(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}], e^i, e^0) \right)_{i \in \llbracket 1, N \rrbracket}, \quad \mathbf{e} = ((e^i)_{i \in \llbracket 1, N \rrbracket}, e^0) \in E^N \times E^0,$$

and reward function

$$\mathbf{f}(\mathbf{x}, \mathbf{a}) = \frac{1}{N} \sum_{i=1}^N f(x^i, a^i, \mu_N[\mathbf{x}, \mathbf{a}]), \quad \mathbf{x} = (x^i)_{i \in \llbracket 1, N \rrbracket}, \quad \mathbf{a} = (a^i)_{i \in \llbracket 1, N \rrbracket}.$$

With respect to standard framework of MDP, we pay a careful attention when dealing with possibly continuous state/action spaces  $(\mathcal{X}, A)$ , and optimizing in general over open-loop controls.

Let us consider the set  $\mathcal{V}$  of sequences  $\boldsymbol{\nu} = (\boldsymbol{\nu}_t)_{t \in \mathbb{N}}$  with  $\boldsymbol{\nu}_0$  a measurable function from  $([0, 1]^N)^{\mathbb{N}}$  into  $A^N$ , and  $\boldsymbol{\nu}_t$  a measurable function from  $([0, 1]^N)^{\mathbb{N}} \times (E^N \times E^0)^t$  into  $A^N$  for  $t \in \mathbb{N}^*$ . For each  $\boldsymbol{\nu} \in \mathcal{V}$ , we can associate a control process  $\boldsymbol{\alpha}^\nu \in \mathcal{A}$  given by

$$\boldsymbol{\alpha}_t^\nu := \boldsymbol{\nu}_t(\mathbf{U}, (\boldsymbol{\varepsilon}_s)_{s \in \llbracket 1, t \rrbracket}), \quad t \in \mathbb{N},$$

(with the convention that  $\boldsymbol{\alpha}_0^\nu = \boldsymbol{\nu}_0(\mathbf{U})$  when  $t = 0$ ), where  $\mathbf{U} = (U_t^i)_{i \in \llbracket 1, N \rrbracket, t \in \mathbb{N}}$  is a family of mutually i.i.d. uniform random variables on  $[0, 1]$ , independent of  $\boldsymbol{\varepsilon}$ , and conversely any control  $\boldsymbol{\alpha} \in \mathcal{A}$  can be represented as  $\boldsymbol{\alpha}^\nu$  for some  $\boldsymbol{\nu} \in \mathcal{V}$ . We call  $\mathcal{V}$  the set of randomized open-loop policies. By misuse of notation, we write  $V_N^\nu = V_N^{\boldsymbol{\alpha}^\nu}$ .

Let us denote by  $L^\infty(\mathcal{X}^N)$  the set of bounded real-valued functions on  $\mathcal{X}^N$ , and by  $L_m^\infty(\mathcal{X}^N)$  the subset of measurable functions in  $L^\infty(\mathcal{X}^N)$ . We then introduce the Bellman ‘‘operator’’  $\mathcal{T}_N : L_m^\infty(\mathcal{X}^N) \rightarrow L^\infty(\mathcal{X}^N)$  defined for any  $W \in L_m^\infty(\mathcal{X}^N)$  by:

$$[\mathcal{T}_N W](\mathbf{x}) := \sup_{\mathbf{a} \in A^N} \mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X}^N.$$

where

$$\mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}) := \mathbf{f}(\mathbf{x}, \mathbf{a}) + \beta \mathbb{E}[W(\mathbf{F}(\mathbf{x}, \mathbf{a}, \boldsymbol{\varepsilon}_1))], \quad \mathbf{x} \in \mathcal{X}^N, \quad \mathbf{a} \in A^N.$$

Notice that the sup can a priori lead to a non measurable function  $\mathcal{T}_N W$ . Because of this,  $\mathcal{T}_N$  is not an operator on  $L_m^\infty(\mathcal{X}^N)$  in the strict sense. To see  $\mathcal{T}_N$  as an operator, we have to find a subset in  $L_m^\infty(\mathcal{X}^N)$  that is preserved by  $\mathcal{T}_N$ . The next result introduces such subset.

**Lemma B.1** *Let  $\mathcal{M}$  be the set in  $L_m^\infty(\mathcal{X}^N)$  defined by*

$$\mathcal{M} := \left\{ W \in L_m^\infty(\mathcal{X}^N) : |W(\mathbf{x}) - W(\mathbf{x}')| \leq 2K_f \sum_{t=0}^{\infty} \beta^t \min[(2K_F)^t \mathbf{d}_N(\mathbf{x}, \mathbf{x}'), \Delta_{\mathcal{X}}], \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}^N \right\}. \quad (\text{B.1})$$

*Then  $\mathcal{M}$  is a complete metric space under the  $\|\cdot\|$  norm, and  $\mathbb{T}_N^{\mathbf{a}}$ , for all  $\mathbf{a} \in A^N$ , and  $\mathcal{T}_N$ , preserve  $\mathcal{M}$ :  $\mathbb{T}_N^{\mathbf{a}} \mathcal{M} \subset \mathcal{M}$ ,  $\mathcal{T}_N \mathcal{M} \subset \mathcal{M}$ .*

**Proof.** It is clear that  $\mathcal{M}$  is closed in  $L_m^\infty(\mathcal{X}^N)$ , and is therefore a complete metric space for  $\|\cdot\|$ . Let  $W \in \mathcal{M}$ . Fix  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^N$ , and  $\mathbf{a} \in A^N$ . Let us start with two preliminary estimations: under  $(\mathbf{Hf}_{\text{lip}})$ , and recalling (2.2), we clearly have

$$|\mathbf{f}(\mathbf{x}, \mathbf{a}) - \mathbf{f}(\mathbf{x}', \mathbf{a})| \leq 2K_f d_N(\mathbf{x}, \mathbf{x}'). \quad (\text{B.2})$$

Similarly, under  $(\mathbf{HF}_{\text{lip}})$ , for  $e^0 \in E^0$ , we have

$$\mathbb{E}[\mathbf{d}_N(\mathbf{F}(\mathbf{x}, \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e^0), \mathbf{F}(\mathbf{x}', \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e^0)] \leq 2K_F \mathbf{d}_N(\mathbf{x}, \mathbf{x}'). \quad (\text{B.3})$$

Thus, denoting by  $\mathbf{X}_1 = \mathbf{F}(\mathbf{x}, \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e^0$  and  $\mathbf{X}'_1 = \mathbf{F}(\mathbf{x}', \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e^0$ , we have, by Jensen's inequality and then (B.3),

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^t \mathbf{d}_N(\mathbf{X}_1, \mathbf{X}'_1), \Delta_{\mathcal{X}}] \right] &\leq \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^t \mathbb{E}[\mathbf{d}_N(\mathbf{X}_1, \mathbf{X}'_1)], \Delta_{\mathcal{X}}] \\ &\leq \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^{t+1} \mathbf{d}_N(\mathbf{x}, \mathbf{x}'), \Delta_{\mathcal{X}}]. \end{aligned} \quad (\text{B.4})$$

The definition of  $\mathbb{T}_N^{\mathbf{a}}$  combined with (B.2), the fact that  $W \in \mathcal{M}$ , and (B.4), implies that

$$\begin{aligned} |\mathbb{T}_N^{\mathbf{a}} W(x) - \mathbb{T}_N^{\mathbf{a}} W(x')| &\leq 2K_f \mathbf{d}_N(\mathbf{x}, \mathbf{x}') + \beta 2K_f \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^{t+1} \mathbf{d}_N(\mathbf{x}, \mathbf{x}'), \Delta_{\mathcal{X}}] \\ &\leq 2K_f \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^t \mathbf{d}_N(\mathbf{x}, \mathbf{x}'), \Delta_{\mathcal{X}}], \end{aligned}$$

which shows that  $\mathbb{T}_N^{\mathbf{a}} W \in \mathcal{M}$ , i.e.  $\mathbb{T}_N^{\mathbf{a}}$  preserves  $\mathcal{M}$ . Furthermore, we have

$$\begin{aligned} |\mathcal{T}_N W(\mathbf{x}) - \mathcal{T}_N W(\mathbf{x}')| &\leq \sup_{\mathbf{a} \in A^N} |\mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}) - \mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}')| \\ &\leq 2K_f \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^t \mathbf{d}_N(\mathbf{x}, \mathbf{x}'), \Delta_{\mathcal{X}}], \end{aligned}$$

which also shows that  $\mathcal{T}_N W \in \mathcal{M}$ . □

Lemma B.1 implies that by restricting  $\mathcal{T}_N$  and  $\mathbb{T}_N^{\mathbf{a}}$  to  $\mathcal{M}$ , we can see  $\mathcal{T}_N$  and  $\mathbb{T}_N^{\mathbf{a}}$  as operators on  $\mathcal{M}$ , that is,  $\mathcal{T}_N : \mathcal{M} \rightarrow \mathcal{M}$  and  $\mathbb{T}_N^{\mathbf{a}} : \mathcal{M} \rightarrow \mathcal{M}$ . However, the property defining the functions in  $\mathcal{M}$  (see (B.1)) is not very natural and practical. The following result provides a more convenient property satisfied by all functions in  $\mathcal{M}$ .

**Lemma B.2** *There exists  $K_\star \in \mathbb{R}$  such that any function  $W \in \mathcal{M}$  is  $\gamma$ -Hölder with constant factor  $K_\star$ , i.e.*

$$|W(\mathbf{x}) - W(\mathbf{x}')| \leq K_\star \mathbf{d}_N(\mathbf{x}, \mathbf{x}')^\gamma, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}^N.$$

**Proof.** We have

$$|W(\mathbf{x}) - W(\mathbf{x}')| \leq 2K_f \sum_{t=0}^{\infty} \beta^t \min [(2K_F)^t \mathbf{d}_N(\mathbf{x}, \mathbf{x}'), \Delta_{\mathcal{X}}] =: 2K_f S(\mathbf{d}_N(\mathbf{x}, \mathbf{x}')).$$

where  $S(m) = \sum_{t=0}^{\infty} \beta^t \min[(2K_F)^t m, \Delta_{\mathcal{X}}]$ . If  $2\beta K_F < 1$ , we clearly have

$$S(m) \leq m \sum_{t=0}^{\infty} (\beta 2K_F)^t = \frac{m}{1 - \beta 2K_F},$$

and so  $W$  is 1-Hölder. Let us now study the case  $2\beta K_F > 1$ . In this case, in particular,  $2K_F > 1$  since  $\beta \in (0, 1)$ , thus  $t \mapsto (2K_F)^t$  is nondecreasing, and so

$$\begin{aligned} S(m) &\leq \sum_{t=0}^{\infty} \int_t^{t+1} \beta^t \min[(2K_F)^s m, \Delta_{\mathcal{X}}] ds \\ &\leq \frac{1}{\beta} \sum_{t=0}^{\infty} \int_t^{t+1} \beta^s \min[(2K_F)^s m, \Delta_{\mathcal{X}}] ds = \frac{1}{\beta} \int_0^{\infty} e^{s \ln \beta} \min[m e^{s \ln(2K_F)}, \Delta_{\mathcal{X}}] ds. \end{aligned}$$

Let  $t_{\star} = t_{\star}(m)$  be such that  $m e^{t_{\star} \ln(2K_F)} = \Delta_{\mathcal{X}}$ , i.e.  $t_{\star} = \frac{\ln(\Delta_{\mathcal{X}}/m)}{\ln(2K_F)}$ . Then,

$$\begin{aligned} \int_0^{\infty} e^{s \ln \beta} \min[m e^{s \ln(2K_F)}, \Delta_{\mathcal{X}}] ds &\leq m \int_0^{t_{\star}} e^{s \ln(2K_F \beta)} ds + \Delta_{\mathcal{X}} \int_{t_{\star}}^{\infty} e^{s \ln(\beta)} ds \\ &= \frac{m}{\ln(2K_F \beta)} \left[ e^{t_{\star} \ln(2K_F \beta)} - 1 \right] - \frac{\Delta_{\mathcal{X}}}{\ln \beta} e^{\ln(\beta) t_{\star}} \\ &= \frac{m}{\ln(2K_F \beta)} \left[ \left( \frac{\Delta_{\mathcal{X}}}{m} \right)^{\frac{\ln(2K_F \beta)}{\ln(2K_F)}} - 1 \right] - \frac{\Delta_{\mathcal{X}}}{\ln \beta} \left( \frac{\Delta_{\mathcal{X}}}{m} \right)^{\frac{\ln(\beta)}{\ln(2K_F)}} \\ &= \Delta_{\mathcal{X}} \left( \frac{1}{\ln(2K_F \beta)} - \frac{1}{\ln \beta} \right) \left( \frac{\Delta_{\mathcal{X}}}{m} \right)^{\frac{\ln(\beta)}{\ln(2K_F)}} - \frac{m}{\ln(2K_F \beta)} \\ &\leq C m^{\min\left[1, \frac{|\ln \beta|}{\ln(2K_F)}\right]} = C m^{\gamma}, \end{aligned}$$

for some positive constant  $C$  depending on  $K_F$ ,  $\beta$  and  $\Delta_{\mathcal{X}}$ . This implies that  $W$  is  $\gamma$ -Hölder with a constant factor  $K_{\star}$  that is clearly independent of  $W \in \mathcal{S}$ . This concludes the proof.  $\square$

The consequence of Lemmas B.1 and B.2 is that the set  $\mathcal{M} \subset L_m^{\infty}(\mathcal{X})$  is a closed set, preserved by  $\mathcal{T}_N$  and contains only functions that are  $\gamma$ -Hölder with factor  $K_{\star}$ . We are now able to get the existence of a unique fixed point to the Bellman operator  $\mathcal{T}_N$ .

**Proposition B.1** (i) The operator  $\mathcal{T}_N$  is monotone increasing: for  $W_1, W_2 \in L_m^{\infty}(\mathcal{X}^N)$ , if  $W_1 \leq W_2$ , then  $\mathcal{T}_N W_1 \leq \mathcal{T}_N W_2$ . (ii) Furthermore, it is contracting on  $L_m^{\infty}(\mathcal{X}^N)$  with Lipschitz factor  $\beta$ , and admits a unique fixed point in  $L_m^{\infty}(\mathcal{X}^N)$ , denoted by  $V_N^*$ , hence solution to:

$$V_N^* = \mathcal{T}_N V_N^*.$$

Moreover,  $V_N^* \in \mathcal{M}$ , and thus  $V_N^*$  is  $\gamma$ -Hölder with constant factor  $K_{\star}$ .

**Proof.** (i) The monotonicity of  $\mathcal{T}_N$  is clear. (ii) The  $\beta$ -contraction property of  $\mathcal{T}_N$  is obtained by standard arguments, which implies the uniqueness of a fixed point (but not the existence). Let us prove the existence of a fixed point. As  $\mathcal{M}$  is preserved by  $\mathcal{T}_N$ , and is closed for  $\|\cdot\|$ , and therefore complete (as a closed subset of the complete space  $L_m^{\infty}(\mathcal{X}^N)$ ), by the Banach fixed point theorem,  $\mathcal{T}_N$  admits a unique fixed point  $V_N^*$  in  $\mathcal{M}$ . By Lemma B.2, this implies that  $V_N^*$  is  $\gamma$ -Hölder with constant factor  $K_{\star}$ , and concludes the proof.  $\square$

**Remark B.1** Notice that the above arguments would not work if we considered, instead of  $\mathcal{M}$ , directly the set of  $\gamma$ -Hölder continuous functions. Indeed, while it is true that such set is stabilized by  $\mathcal{T}_N$  (it essentially follows from (B.2) and (B.3)), the set of  $\gamma$ -Hölder continuous functions is not closed in  $L_m^\infty(\mathcal{X}^N)$  (and thus not a complete metric space): there might indeed exist a converging sequence of  $\gamma$ -Hölder continuous functions with multiplicative factors (in the Hölder property) tending toward infinity, such that the limit function is not  $\gamma$ -Hölder anymore.  $\square$

As a consequence of Proposition B.1, we can show the following relation between the value function  $V_N$  of the  $N$ -agent MDP, and the fixed point  $V_N^*$  of the Bellman operator  $\mathcal{T}_N$ .

**Lemma B.3** *For all  $\mathbf{x} \in \mathcal{X}^N$ , we have  $V_N(\mathbf{x}) \leq V_N^*(\mathbf{x})$ .*

**Proof.** For any  $\mathbf{x} \in \mathcal{X}^N$ ,  $\boldsymbol{\nu} \in \mathcal{V}$ , we have

$$\begin{aligned} \mathbb{E}\left[\mathbf{f}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{U})) + \beta V_N^*(\mathbf{F}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{U}), \boldsymbol{\varepsilon}_1))\right] &= \mathbb{E}\left[\left\{\mathbf{f}(\mathbf{x}, \boldsymbol{\nu}_0(u)) + \beta \mathbb{E}[V_N^*(\mathbf{F}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{u}), \boldsymbol{\varepsilon}_1))]\right\}_{\mathbf{u}:=\mathbf{U}}\right] \\ &= \mathbb{E}\left[\mathbb{T}^{\boldsymbol{\nu}_0} V_N^*(\mathbf{x})\right] \leq \mathcal{T} V_N^*(\mathbf{x}) = V_N^*(\mathbf{x}). \end{aligned} \quad (\text{B.5})$$

For any  $(\mathbf{u}, \mathbf{e}) \in ([0, 1]^N)^{\mathbb{N}} \times (E^N \times E^0)$ , and for any  $\boldsymbol{\nu} \in \mathcal{V}$ , we define  $\vec{\boldsymbol{\nu}}^{\mathbf{u}, \mathbf{e}} \in \mathcal{V}$  by

$$\vec{\boldsymbol{\nu}}_t^{\mathbf{u}, \mathbf{e}}(\mathbf{u}', (\mathbf{e}'_s)_{s \in [1, t]}) := \boldsymbol{\nu}_{t+1}(\mathbf{u}, \mathbf{e}, (\mathbf{e}'_s)_{s \in [1, t]}), \quad (\mathbf{u}', (\mathbf{e}'_s)_{s \in [1, t]}) \in ([0, 1]^N)^{\mathbb{N}} \times (E^N \times E^0)^t, t \in \mathbb{N}.$$

Standard Markov arguments imply the following flow property for randomized open-loop policies:

$$V_N^\boldsymbol{\nu}(\mathbf{x}) = \mathbb{E}\left[\mathbf{f}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{U})) + \beta V_N^{\vec{\boldsymbol{\nu}}^{\mathbf{U}, \boldsymbol{\varepsilon}_1}}(\mathbf{F}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{U}), \boldsymbol{\varepsilon}))\right].$$

Together with (B.5), we then get

$$\begin{aligned} V_N^*(\mathbf{x}) - V_N^\boldsymbol{\nu}(\mathbf{x}) &\geq \beta \mathbb{E}\left[V_N^*(\mathbf{F}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{U}), \boldsymbol{\varepsilon}_1) - V_N^{\vec{\boldsymbol{\nu}}^{\mathbf{U}, \boldsymbol{\varepsilon}_1}}(\mathbf{F}(\mathbf{x}, \boldsymbol{\nu}_0(\mathbf{U}), \boldsymbol{\varepsilon}_1))\right] \\ &\geq \beta \inf_{\mathbf{x} \in \mathcal{X}^N, \boldsymbol{\nu} \in \mathcal{V}} \{V_N^*(\mathbf{x}) - V_N^\boldsymbol{\nu}(\mathbf{x})\}. \end{aligned}$$

Taking the infimum over  $\mathbf{x} \in \mathcal{X}^N$ ,  $\boldsymbol{\nu} \in \mathcal{V}$  on the left hand side of the above inequality, and since  $\beta < 1$ , this shows that  $V_N^\boldsymbol{\nu}(\mathbf{x}) \leq V_N^*(\mathbf{x})$  for all  $\boldsymbol{\nu} \in \mathcal{V}$ . We conclude that  $V_N \leq V_N^*$ .  $\square$

We aim now to prove rigorously the equality  $V_N = V_N^*$ , i.e., the value function  $V_N$  of the  $N$ -agent MDP satisfies the Bellman fixed point equation:  $V_N = \mathcal{T}_N V_N$ , and also to show the existence of  $\varepsilon$ -optimal (randomized) feedback control for  $V_N$ .

A feedback policy (resp. randomized feedback policy) is an element  $\boldsymbol{\pi} \in L^0(\mathcal{X}^N; A^N)$  (resp.  $L^0(\mathcal{X}^N \times [0, 1]^N; A^N)$ ), the set of measurable functions from  $\mathcal{X}^N$  (resp.  $\mathcal{X}^N \times [0, 1]^N$ ) into  $A^N$ . The associated feedback control is the unique control  $\boldsymbol{\alpha}^\boldsymbol{\pi}$  given by  $\boldsymbol{\alpha}_t^\boldsymbol{\pi} = \boldsymbol{\pi}(\mathbf{X}_t)$ , (resp.  $\boldsymbol{\pi}_\tau(\mathbf{X}_t, \mathbf{U}_t)$ ),  $t \in \mathbb{N}$ , where  $\{\mathbf{U}_t = (U_t^i)_{i \in [1, N]}, t \in \mathbb{N}\}$  is a family of mutually i.i.d. uniform random variables on  $[0, 1]$ , independent of  $\mathcal{G}$ ,  $\boldsymbol{\varepsilon}$ . By misuse of notation, we denote  $V_N^\boldsymbol{\pi} = V_N^{\boldsymbol{\alpha}^\boldsymbol{\pi}}$ . Given  $\boldsymbol{\pi} \in L^0(\mathcal{X}^N; A^N)$  (resp.  $L^0(\mathcal{X}^N \times [0, 1]^N; A^N)$ ), we introduce the operator  $\mathcal{T}_N^\boldsymbol{\pi}$  on  $L_m^\infty(\mathcal{X}^N)$ , defined by

$$\mathcal{T}_N^\boldsymbol{\pi} W(\mathbf{x}) := \mathbf{f}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \beta \mathbb{E}[W(\mathbf{F}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}), \boldsymbol{\varepsilon}_1))], \quad \mathbf{x} \in \mathcal{X}^N,$$



resp.

$$\mathcal{T}_N^\pi W(\mathbf{x}) := \mathbb{E}[\mathbf{f}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}, \mathbf{U}_0)) + \beta W(\mathbf{F}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}, \mathbf{U}_0), \boldsymbol{\varepsilon}_1))], \quad \forall \mathbf{x} \in \mathcal{X}^N,$$

where  $\mathbf{U}_0 = (U_0^i)_{i \in [1, N]}$  is a family of i.i.d.  $\sim \mathcal{U}([0, 1])$ , independent of  $\mathcal{G}, \boldsymbol{\varepsilon}$ .

We have the basic and standard properties on the operator  $\mathcal{T}_N^\pi$ :

**Lemma B.4** Fix  $\boldsymbol{\pi} \in L^0(\mathcal{X}^N; A^N)$  (resp.  $L^0(\mathcal{X}^N \times [0, 1]^N; A^N)$ ).

- (i) The operator  $\mathcal{T}_N^\pi$  is  $\beta$ -contracting on  $L_m^\infty(\mathcal{X}^N)$ , and  $V_N^\pi$  is its unique fixed point.
- (ii) Furthermore, it is monotone increasing: for  $W_1, W_2 \in L^\infty(\mathcal{X}^N)$ , if  $W_1 \leq W_2$ , then  $\mathcal{T}_N^\pi W_1 \leq \mathcal{T}_N^\pi W_2$ .

We state the standard verification type result for the  $N$ -individual MDP, by means of the Bellman operator.

**Lemma B.5 (Verification result)**

Fix  $\epsilon \geq 0$ , and suppose that there exists an  $\epsilon$ -optimal (randomized) feedback policy  $\boldsymbol{\pi}^\epsilon$  for  $V_N^*$  in the sense that

$$V_N^* \leq \mathcal{T}_N^{\boldsymbol{\pi}^\epsilon} V_N^* + \epsilon.$$

Then,  $\boldsymbol{\alpha}^{\boldsymbol{\pi}^\epsilon} \in \mathcal{A}$  is  $\frac{\epsilon}{1-\beta}$ -optimal for  $V_N$ , i.e.,  $V_N^{\boldsymbol{\pi}^\epsilon} \geq V_N - \frac{\epsilon}{1-\beta}$ , and we have  $V_N \geq V_N^* - \frac{\epsilon}{1-\beta}$ .

**Proof.** Since  $V_N^{\boldsymbol{\pi}^\epsilon} = \mathcal{T}_N^{\boldsymbol{\pi}^\epsilon} V_N^{\boldsymbol{\pi}^\epsilon}$ , and recalling from Lemma B.3 that  $V_N^* \geq V_N \geq V_N^{\boldsymbol{\pi}^\epsilon}$ , we have for all  $\mathbf{x} \in \mathcal{X}^N$ ,

$$\left| (V_N^* - V_N^{\boldsymbol{\pi}^\epsilon})(\mathbf{x}) \right| \leq \left| \mathcal{T}_N^{\boldsymbol{\pi}^\epsilon} V_N^*(\mathbf{x}) - \mathcal{T}_N^{\boldsymbol{\pi}^\epsilon} V_N^{\boldsymbol{\pi}^\epsilon}(\mathbf{x}) \right| + \epsilon \leq \beta \|V_N^* - V_N^{\boldsymbol{\pi}^\epsilon}\| + \epsilon,$$

where we used the  $\beta$ -contraction property of  $\mathcal{T}_N^{\boldsymbol{\pi}^\epsilon}$  in Lemma B.4. We deduce that  $\|V_N^* - V_N^{\boldsymbol{\pi}^\epsilon}\| \leq \frac{\epsilon}{1-\beta}$ , and then,  $V_N \geq V_N^{\boldsymbol{\pi}^\epsilon} \geq V_N^* - \frac{\epsilon}{1-\beta}$ , which combined with  $V_N^* \geq V_N$ , concludes the proof.  $\square$

We finally conclude this section by showing the existence of an  $\epsilon$ -optimal (randomized) feedback policy for  $N$ -agent MDP on  $\mathcal{X}^N$ , and obtain as a by-product the corresponding Bellman fixed point equation for its value function.

**Proposition B.2** For all  $\epsilon > 0$ , there exists a (randomized) feedback policy  $\boldsymbol{\pi}^\epsilon$  that is  $\epsilon$ -optimal for  $V_N^*$ . Consequently, the control  $\boldsymbol{\alpha}^{\boldsymbol{\pi}^\epsilon} \in \mathcal{A}$  is  $\frac{\epsilon}{1-\beta}$ -optimal for  $V_N$ , and we have  $V_N = V_N^*$ , which thus satisfies the Bellman fixed point equation.

**Proof.** We prove the result for  $\epsilon$ -optimal feedback policy (the case of  $\epsilon$ -optimal randomized feedback policy is dealt with similarly). Fix  $\epsilon > 0$ , and given  $\eta > 0$ , consider a quantizing grid  $\mathcal{M}^\eta = \{\mathbf{x}_1, \dots, \mathbf{x}_{N^\eta}\} \subset \mathcal{X}^N$ , and an associated partition  $C_k^\eta, k = 1, \dots, N^\eta$ , of  $\mathcal{X}^N$ , satisfying

$$C_k^\eta \subset B^\eta(\mathbf{x}_k) := \left\{ \mathbf{x} \in \mathcal{X}^N : \mathbf{d}_N(\mathbf{x}, \mathbf{x}_k) \leq \eta \right\}, \quad k = 1, \dots, N^\eta.$$

For any  $\mathbf{x}_k, k = 1, \dots, N^\eta$ , there exists  $\mathbf{a}_k^\epsilon \in A^N$  such that

$$V_N^*(\mathbf{x}_k) \leq \mathbb{T}^{\mathbf{a}_k^\epsilon} V_N^*(\mathbf{x}_k) + \frac{\epsilon}{3}. \tag{B.6}$$

From the partition  $C_k^\eta$ ,  $k = 1, \dots, N_\eta$  of  $\mathcal{X}^N$ , associated to  $\mathcal{M}^\eta$ , we construct the function  $\pi^\epsilon : \mathcal{X}^N \rightarrow A^N$  as follows: we define, for all  $\mathbf{x} \in \mathcal{X}^N$ ,

$$\pi^\epsilon(\mathbf{x}) = \mathbf{a}_k^\epsilon, \quad \text{when } \mathbf{x} \in C_k^\eta, \quad k = 1, \dots, N_\eta.$$

Such function  $\pi^\epsilon$  is clearly measurable. Let us now check that such  $\pi^\epsilon$  yields an  $\epsilon$ -optimal feedback policy for  $\eta$  small enough. For  $\mathbf{x} \in \mathcal{X}^N$ , we define  $\mathbf{x}^\eta = \mathbf{x}_k$ , when  $\mathbf{x} \in C_k^\eta$ ,  $k = 1, \dots, N_\eta$ . Observe that  $d_N(\mathbf{x}, \mathbf{x}^\eta) \leq \eta$ . We then write for any  $\mathbf{x} \in \mathcal{X}^N$ ,

$$\begin{aligned} [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) - V_N^*(\mathbf{x}) &= \left( [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) - [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}^\eta) \right) + \left( [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}^\eta) - V_N^*(\mathbf{x}^\eta) \right) \\ &\quad + (V_N^*(\mathbf{x}^\eta) - V_N^*(\mathbf{x})) \\ &\geq \left( [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) - [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}^\eta) \right) - \frac{\epsilon}{3} - \frac{\epsilon}{3}, \end{aligned} \quad (\text{B.7})$$

where we used (B.6) and the fact that  $|V_N^*(\mathbf{x}^\eta) - V_N^*(\mathbf{x})| \leq \epsilon/3$  for  $\eta$  small enough by uniform continuity of  $V_N^*$  in Proposition B.1. Moreover, by observing that  $\pi^\epsilon(\mathbf{x}) = \pi^\epsilon(\mathbf{x}^\eta) =: \mathbf{a}$ , we have

$$\begin{aligned} [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) &= \mathbb{E} \left[ \mathbf{f}(\mathbf{x}, \mathbf{a}) + \beta V_N^*(\mathbf{F}(\mathbf{x}, \mathbf{a}, \varepsilon_1)) \right], \\ [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}^\eta) &= \mathbb{E} \left[ \mathbf{f}(\mathbf{x}^\eta, \mathbf{a}) + \beta V_N^*(\mathbf{F}(\mathbf{x}^\eta, \mathbf{a}, \varepsilon_1)) \right]. \end{aligned}$$

Under  $(\mathbf{HF}_{\text{lip}})$ - $(\mathbf{Hf}_{\text{lip}})$ , and by using the  $\gamma$ -Hölder property of  $V_N^*$  with constant  $K_\star$  in Proposition B.1, we then get

$$\begin{aligned} &|[\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) - [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}^\eta)| \\ &\leq 2K_f d_N(\mathbf{x}, \mathbf{x}^\eta) + \beta K_\star \mathbb{E} \left[ \mathbb{E} \left[ d_N(\mathbf{F}(\mathbf{x}, \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e), \mathbf{F}(\mathbf{x}^\eta, \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e) \right]^\gamma \right]_{e:=\varepsilon_1^0} \\ &\leq 2K_f d_N(\mathbf{x}, \mathbf{x}^\eta) + \beta K_\star \mathbb{E} \left[ \mathbb{E} \left[ d_N(\mathbf{F}(\mathbf{x}, \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e), \mathbf{F}(\mathbf{x}^\eta, \mathbf{a}, (\varepsilon_1^i)_{i \in [1, N]}), e) \right]^\gamma \right]_{e:=\varepsilon_1^0} \\ &\leq C d_N(\mathbf{x}, \mathbf{x}^\eta)^\gamma \leq C \eta^\gamma. \end{aligned}$$

for some constant  $C$ . Therefore,  $|[\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) - [\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}^\eta)| \leq \epsilon/3$ , and, plugging into (B.7), we obtain  $[\mathcal{T}_N^{\pi^\epsilon} V_N^*](\mathbf{x}) - V_N^*(\mathbf{x}) \geq -\epsilon$ , for all  $\mathbf{x} \in \mathcal{X}^N$ , which means that  $\pi^\epsilon$  is  $\epsilon$ -optimal for  $V_N^*$ . The rest of the assertions in the Theorem follows from the verification result in Lemma B.5.  $\square$

## References

- [1] N. Bäuerle. Mean Field Markov Decision Processes. *arXiv preprint arXiv:2106.08755*, to appear in *Applied Mathematics and Optimization*, 2021.
- [2] D. Bertsekas and S. E Shreve. *Stochastic optimal control: the discrete-time case*, volume 5. Athena Scientific, 1996.
- [3] E. Boissard and T. Le Gouic. On the mean speed of convergence of empirical and occupation measures in Wasserstein distance. In *Annales de l'IHP Probabilités et statistiques*, volume 50, pages 539–563, 2014.
- [4] P. Cardaliaguet, S. Daudin, J. Jackson, and P. E. Souganidis. An algebraic convergence rate for the optimal control of McKean-Vlasov dynamics. *arXiv:2203.14554*, 2022.

- [5] R. Carmona and F. Delarue. *Probabilistic Theory of Mean Field Games with Applications, vol I and II*. Probability theory and stochastic modelling. Springer, 2018.
- [6] R. Carmona, M. Laurière, and Z. Tan. Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. *arXiv preprint arXiv:1910.12802*, 2019.
- [7] A. Cecchin. Finite state  $N$ -agent and mean field control problems. *ESAIM: Control, Optimization and Calculus of Variations*, 27(31), 2021.
- [8] M. F. Djete. Extended mean field control problem: a propagation of chaos result. *arXiv preprint arXiv:2006.12996*, to appear in *Electronic Journal of Probability*, 2020.
- [9] N. Fournier and A. Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3):707–738, 2015.
- [10] W. Gangbo, S. Mayorga, and A. Swiech. Finite dimensional approximations of Hamilton-Jacobi-Bellman equations in spaces of probability measures. *SIAM. J. Math. Anal*, 2:1320–1356, 2021.
- [11] N. Gast and B. Gaujal. A mean field approach for optimization in discrete time. *Discrete Event dynamics Systems*, 21(1):63–101, 2011.
- [12] M. Germain, H. Pham, and X. Warin. Rate of convergence for particles approximation of PDEs in Wasserstein space. *arXiv:2103.00837*, to appear in *Journal of Applied Probability*, 2021.
- [13] O. Kallenberg. *Random measures, theory and applications*, volume 1. Springer, 2017.
- [14] D. Lacker. Limit theory for controlled McKean-Vlasov dynamics. *SIAM Journal on Control and Optimization*, 55(3):1641–1672, 2017.
- [15] M. Motte. *Mathematical models for large populations, behavioral economics, and targeted advertising*. PhD thesis, Université Paris Cité, 2021.
- [16] M. Motte and H. Pham. Mean-field Markov decision processes with common noise and open-loop controls. *arXiv:1912.07883*, to appear in *Annals of Applied Probability*, 2019.
- [17] H. Pham and X. Wei. Discrete time McKean-Vlasov control problem: a dynamic programming approach. *Applied Mathematics & Optimization*, 74(3):487–506, 2016.
- [18] A.S. Sznitman. Topics in propagation of chaos. In Springer, editor, *Ecole d’été de probabilités de Saint Flour, lecture notes in Mathematics*, volume 1464, pages 165–251. 1989.
- [19] M. Thorpe. Introduction to optimal transport, 2018.