

1 **Response about the significance and originality.** Modelling the dynamics of multi-agent learning has long been  
2 an important research topic, but an  $n$ -agent setting where  $n$  tends to infinity has not been considered. All of the  
3 previous works focus on 2-agent settings and mostly use evolutionary game theoretic approaches (see the recent survey  
4 [Bloembergen et al., JAIR'16]). Our mean field theoretic approach is of fundamental difference from the evolutionary  
5 game theoretic approaches. As we explain in Introduction, the use of evolutionary game theory is inappropriate for  
6  $n$ -agent settings, because, in principle, the number of equations required to model the entire population dynamics is  
7 proportional to the number of agents in the population. As  $n$  tends to infinity, analyzing or solving this system of  
8 equations becomes practically infeasible. In this paper, we show that by using mean field theory, only three equations are  
9 required to describe the dynamics of the whole population. A system of such small number of equations, as presented  
10 in Eq. 17, greatly reduces the problem complexity and makes the modelling tractable. Therefore, this paper introduces  
11 a new methodology to the modelling of learning dynamics in an infinitely large agent population, which is an emerging  
12 research topic given the growing interest in large-scale multi-agent systems.

13 The theoretical contributions of our work and the works [Mguni et al., AAAI'18; Mguni et al., AAMAS'19] mentioned  
14 by Reviewer 1 are very different. The works of Mguni et al. propose novel learning or incentive design methods, and  
15 prove that these methods will finally result in the convergence to efficient Nash equilibria in an infinitely large agent  
16 population. The actual process of convergence, however, is not formally described. This paper, to the best of our  
17 knowledge, is the first time to formally show the reinforcement learning dynamics, say, how the policies of individual  
18 agents gradually evolve over time, in an infinitely large agent population. In particular, the heart of this paper – a  
19 Fokker-Planck equation describing the evolution of the probability distribution of  $Q$ -values in an agent population – has  
20 not been reported elsewhere.

21 In this paper, we focus our attention on the population dynamics of an infinitely large agent population that use  
22  $Q$ -learning. This is because  $Q$ -learning is one of the most important algorithms in reinforcement learning research and  
23 is the basis of a number of multi-agent reinforcement learning algorithms. Considering other learning algorithms will  
24 be an interesting and also plausible extension to our work.

25 We apologize that the above points should have been clearer. We shall highlight these points in the revised version.

26 **Response about the experiments.** The experimental study of this paper aims to illustrate and validate our mean  
27 field theoretic model. The games we select (prisoner's dilemma, stag hunt, hawk dove and choosing side) are typical  
28 matrix games that vary in the number, symmetry and efficiency of Nash equilibrium. This makes them good examples  
29 for illustration and validation, since they can be easily understood, but will lead to qualitatively different patterns of  
30 population dynamics. The nearly precise matching of the population dynamics derived from our model to those obtained  
31 from the agent-based simulations for each game type provides a clear and effective validation of our model. To further  
32 exhibit the strength of our model, we will find more complicated yet still understandable games to experiment on in the  
33 revised version. We will release our codes of the experiments if this paper get accepted.

34 **Response to the questions raised by Reviewer 2.** In each entry of Table 1, the first number is the payoff of the row  
35 player, while the second one is that of the column player. We shall include this specification in the revised version. The  
36 term  $\gamma \max_{a' \in A} Q_t^{s', a'}(n_i)$  in Eq. 1 estimates the optimal discounted future payoff of player  $n_i$  under state  $s'$ , after it  
37 plays action  $a$  under the current state  $s$  and consequently transits to the new state  $s'$ . For a matrix game, at a given time  
38 step  $t$ , agents play one round of the game. The row and the column players each takes one action simultaneously and  
39 receives an immediate payoff based on the joint actions. Then, the game ends. At the next time step  $t + 1$ , the agents  
40 play a new round of the game. In other words, from time  $t$  to  $t + 1$ , there is no state transition for an agent. Since  
41 there is no state transition at all, there is no need to maintain the term  $\gamma \max_{a' \in A} Q_t^{s', a'}(n_i)$ . Hence, it is a common  
42 practice to remove this term from the  $Q$ -value update function for matrix games [Gomes and Kowalczyk, ICML'09;  
43 Wunder et al., ICML'10; Kianercy et al., Physical Review E'12]. We apologize for the lack of explanation in the current  
44 version, and shall provide a detailed one in the revised version. We appreciate a lot for pointing out our typos and giving  
45 valuable suggestions on the language!

46 **Response to the questions raised by Reviewer 3.** The agent-based simulations are conducted on 100 agent populations  
47 each consists of 1,000 agents. Agents play games strictly following the interaction scheme presented in Algorithm 1,  
48 and use  $Q$ -learning to update their policies. In Eq. 4 and 5,  $a$  should have been  $a_j$ . Eq. 5 holds for any valid value of  $\eta$ ,  
49 which, by the definition of learning rate should be between 0 and 1. In Eq. 8, the series should be convergent, since  
50 the function  $u(a, \mathbf{x}_t(n_j))$  is an analytic function. Given each element of the vector  $\Delta \mathbf{x}_t(n_j)$  is between 0 and 1, we  
51 consider the second order and the higher order terms negligible. When  $m, n \rightarrow \infty$ , Eq. 8 holds. From Eq. 9, we can  
52 tell that the trajectory of each agent depends on its  $Q$ -values and is independent of who he/she is. Hence, we consider  
53 the trajectories of agents to be a function of  $Q$ -values in Eq. 10. We shall rewrite the left hand side of Eq. 10 to be  
54  $\mathbb{E}[\frac{dQ_t^{a_i}}{dt}](\mathbf{q}_t)$  for clarity. Given  $n \rightarrow \infty$ , the state of the population can be characterized by a distribution of  $Q$ -values in  
55 the population. Therefore, by deriving the Fokker-Planck equation that describes the time evolution of the  $Q$ -value  
56 distribution, we show in Eq.17 that only three equations are required to describe the entire population dynamics.