

# A simulation framework for measuring robustness of incentive mechanisms and its implementation in reputation systems

Yuan Liu<sup>1</sup> · Jie Zhang<sup>1</sup> · Bo An<sup>1</sup> · Sandip Sen<sup>2</sup>

Published online: 3 April 2015  
© The Author(s) 2015

**Abstract** In game theoretical analysis of incentive mechanisms, all players are assumed to be rational. Since it is likely that mechanism participants in the real world may not be fully rational, such mechanisms may not work as effectively as in the idealized settings for which they were designed. Therefore, it is important to evaluate the robustness of incentive mechanisms against various types of agents with bounded rational behaviors. Such evaluations would provide us with the information needed to choose mechanisms with desired properties in real environments. In this article, we first propose a general robustness measure, inspired by research in evolutionary game theory, as the maximal percentage of invaders taking non-equilibrium strategies such that the agents sustain the desired equilibrium strategy. We then propose a simulation framework based on evolutionary dynamics to empirically evaluate the equilibrium robustness. The proposed simulation framework is validated by comparing the simulated results with the analytical predictions based on a modified simplex analysis approach. Finally, we implement the proposed simulation framework for evaluating the robustness of incentive mechanisms in reputation systems for electronic marketplaces. The results from the implementation show that the evaluated mechanisms have high robustness against a certain non-equilibrium strategy, but is vulnerable to another strategy, indicating the need for designing more robust incentive mechanisms for reputation management in e-marketplaces.

---

✉ Yuan Liu  
liu.yuan@ntu.edu.sg

Jie Zhang  
zhangj@ntu.edu.sg

Bo An  
boan@ntu.edu.sg

Sandip Sen  
sandip-sen@utulsa.edu

<sup>1</sup> School of Computer Engineering, Nanyang Technological University, Singapore, Singapore

<sup>2</sup> Department of Mathematical and Computer Sciences, The University of Tulsa, Tulsa, USA

**Keywords** Robustness · Irrational agents · Incentive mechanism · Reputation system · Electronic marketplace

## 1 Introduction

Game theory provides a powerful theoretical framework to analyze strategic interactions where the payoff of a player depends on his/her strategy and those of others. Incentive mechanisms, based on game theoretical analysis, have been designed to promote the desired behavior of rational players. For example, consider reputation systems in electronic commerce applications, where the ratings provided by independent buyers are aggregated to assist other buyers in choosing satisfactory sellers. The reputation system desires truthful ratings from buyers, and incentive mechanisms, such as the side-payment incentive mechanism [13], have been proposed where the equilibrium strategy of rational buyers is to provide truthful ratings. However, in realistic scenarios, players with bounded rationality [29] may deviate from the desired equilibrium strategy, causing other rational players to also deviate, and as a result the incentive mechanisms may fail to achieve the expected performance.

The study of bounded rationality is increasingly important in the field of mechanism design. The existence of bounded rational players, who may make mistakes due to reasons unknown to mechanism designers, has been widely recognized by many researchers [3, 5, 35]. Aiyer et al. [2] propose a filtering approach to eliminate the impact of naive Byzantine (irrational) players who randomly deviate from equilibrium strategies. Moreover, the qualitative analysis of the robustness against various attacks has been proposed, such as [34], where a mechanism is claimed to be either *resistant* or *vulnerable* with respect to an attacking strategy. However, one common limitation of the existing approaches is that they cannot tell to what extent these incentive mechanisms are robust against non-equilibrium strategies. Measuring the robustness of mechanisms is useful for ensuring the practical usability of incentive mechanisms and towards the design of more robust incentive mechanisms.

In this article, we aim to propose a qualitative approach to evaluate the robustness of incentive mechanisms against bounded rational players who implement various non-equilibrium strategies. We first define a general robustness measure of a desired equilibrium, inspired by evolutionary game theory, as the maximum percentage of bounded rational agents existing in the system while it is still better off for rational agents to perform the suggested strategies in the desired equilibrium. A practical robustness measure with feasible computational cost is then proposed by considering a finite number of rational agents and a specific non-equilibrium strategy. Due to the complex settings of various mechanisms and various non-equilibrium strategies existing in a realistic case, a simulation framework is proposed to calculate the proposed robustness. We then formalize various forms of representative games in a uniform format and validate the proposed simulation framework. It is shown to be able to produce the same results as those by analytical analysis. Finally, we implement our simulation framework in evaluating and comparing the robustness of four incentive mechanisms for reputation systems in e-marketplaces where bounded rational agents adopt different attacking strategies. We have observed that the robustness of mechanisms decreases against sophisticated non-equilibrium strategies, demonstrating the need for more robust incentive mechanisms for reputation systems.

The key contributions of this article include:

1. a general robustness measure of a desired equilibrium and a practical measurement against a specific non-equilibrium strategy;

2. a simulation framework for quantitatively evaluating the practical robustness measurement;
3. the validation of the simulation framework in representative games;
4. the implementation of the framework in evaluating incentive mechanisms in reputation systems for e-marketplaces against various types of untruthful attacking strategies.

The remainder of this article is structured as follows. A review of related work is given in Sect. 2. In Sect. 3, two formal robustness measures for incentive mechanisms are defined when the populations are infinite and finite, respectively. In Sect. 4, a simulation framework for evaluating the robustness is proposed. The validation of the simulation framework in representative games is presented in Sect. 5.2, and the implementation and evaluation of the simulation framework in e-marketplaces are presented in Sect. 6. Finally, Sect. 7 concludes the article with an overview of future work.

## 2 Related work

In classical game theory, the players or agents are perfectly rational and behave according to their and other's preferences or payoffs in their interactions [23,30]. Players in the classical setting generally have a perfect knowledge of the environment and the payoffs and try to maximise their individual utility. However, under the biological circumstances, e.g. a population of birds or insects surviving on an isolated island, it becomes impossible to judge what actions are rational. Instead, players learn to optimise their behaviors and maximise their return through a process of trial and error [27]. This learning process matches the concept of evolution in biology, and forms the basis of evolutionary game theory [23,24,32]. Based on this, in this article, we aim to propose a robustness metric where some agents are not rational.

In the literature of incentive mechanism design, increasing attention has been paid to designing practical incentive mechanisms which could achieve the desired outcome while relaxing the assumption that every player is rational (e.g., [1,3,6]). These studies aim to design mechanisms for players with bounded rationality. However, there are various reasons for a player being bounded rational, e.g., unknown external utility or emotion [7,29], and it is nearly impossible to motivate them to behave "rationally". We adopt a more practical perspective and focus on studying the behavior of rational players when bounded rational players choose non-equilibrium strategies.

The robustness problem is also referred to as the implementation problem where some of the players are "faulty" in the sense that they fail to act optimally [5]. The term  $k$ -FTNE is defined to serve as a stronger equilibrium concept than the Nash equilibrium, where each non-faulty player would sustain the designed strategy as long as other non-faulty players sustain the designed strategy, regardless of the strategy used by up to  $k$  faulty players. Another relevant concept in game theory is  $k-t$  robust Nash equilibrium (NE) [6]. An equilibrium is  $t$ -immune if it can be tolerant up to  $t$  irrational agents. An equilibrium is  $k$ -resilient if all the rational agents in a coalition of size up to  $k$  ( $k \geq 1$ ) will not deviate from the equilibrium strategy. The limitation of the above studies is that they only provide an abstract characterization of a robust equilibrium, and does not provide an operational approach to quantitatively evaluating mechanism robustness. In this article, we formally define the measurement of robustness for incentive mechanisms and propose a simulation framework to quantitatively evaluate the robustness.

To improve the robustness and protect incentive mechanisms from the influence of bounded rational players, a crowd filtering approach has been proposed. In [2], a fault tolerance model

named BAR model is used to filter out players who randomly deviate from the desired equilibrium strategies. However, the model fails to filter out players with complex strategies. Since it is difficult to insulate the impact of every non-equilibrium strategy, it is necessary to study the robustness of mechanisms and our work meets that need.

Qualitative analysis of the robustness of incentive mechanisms has been proposed for different non-equilibrium strategies. Witkowski et al. [34] analyze the conditions when a mechanism is robust against various types of collusive attacks. In contrast, we aim to provide a quantitative measurement of robustness, which is especially useful for comparing incentive mechanisms used in practice.

In reputation systems for e-marketplaces, many incentive mechanisms have been widely used to promote truthful feedbacks (ratings) from agents, such as side-payment mechanism [13], credibility mechanism [19], and trust-based incentive mechanism [36]. They are designed based on different methodology, and have been proven that the strategy of truthful reporting is an equilibrium strategy for each agent. However, there is no measurement to justify which one is more practical when they are implemented in a realistic environment. In this article, we implement the proposed simulation framework to evaluate and compare the robustness of these mechanisms against various non-equilibrium strategies. We can not find a dominant mechanism among these, as for every pair of mechanisms there exist different scenarios where each member of the pair is more robust. This result establishes the need for a more robust incentive mechanism and our work provides an evaluation measurement for researchers working on developing such a mechanism.

### 3 Robustness measure

Incentive mechanism design is a sub-field of game theory that considers how to implement the desired, e.g., socially efficient, solutions by motivating rational players to truthfully disclose their private information [10, 15]. Players' equilibrium strategy profile<sup>1</sup> is called a desired strategy profile (equilibrium) which leads to the desired outcome of the mechanism, assuming that the players are rational. Players who instead adopt a non-equilibrium strategy are recognized to be *bounded rational*. We study the robustness of incentive mechanisms when only some of the players are bounded rational.

Our robustness measure of incentive mechanisms is inspired by the concept of evolutionary stable strategy (ESS) studied in evolutionary game theory—the application of game theory to biology [21, 22, 25, 32]—to study the evolution of strategies in the population(s) of individuals who are competing with each other for survival and reproduction. A strategy is evolutionarily stable if a population of individuals using the strategy will not deviate from it even when the population is invaded or infiltrated by a small fraction of mutants, i.e., individuals using an alternative (non-equilibrium) strategy. In an incentive mechanism, if the desired equilibrium strategy is an ESS, the mechanism will be robust against a sufficiently small fraction of bounded rational invaders. Furthermore, an incentive mechanism is more robust than another if its desired equilibrium is robust against a larger fraction of invaders. Based on this intuition, we propose a general quantitative definition of robustness with respect to the desired equilibrium.

We work within an evolutionary framework of agent populations. Within this framework an incentive mechanism in a game with  $n$  players  $I = \{1, \dots, n\}$  is modeled by  $n$  popula-

---

<sup>1</sup> For ease of analysis, we assume a single equilibrium for a mechanism. Our analysis can be extended to handle multiple equilibria.

tions: for each player position,<sup>2</sup> there is a large population of individuals (agents), and each such individual is following a pure strategy. Let  $S_i = \{s_1^i, \dots, s_{m_i}^i\}$  be the set of pure strategies available to the individuals of population  $i$  and  $\Delta_i = \{x_i \in \mathbb{R}^{m_i} \mid \sum_{s_j^i \in S_i} x_i(j) = 1, x_i(j) \geq 0, j = 1, \dots, m_i\}$  be the set of possible strategy profiles for population  $i$ , where each  $x_i(j)$  corresponds to the fraction of individuals in the population  $i$  playing strategy  $s_j^i \in S_i$ . Note that  $x_i$  is formally equivalent to a mixed strategy for the player  $i \in I$  in the  $n$ -player game. The combination of  $n$  population profiles is  $x \in \Theta$  where  $\Theta = \times_{i \in I} \Delta_i$ . Suppose  $\epsilon \in (0, 1)$  proportion of the population  $i$  invade the mechanism by playing a different profile  $y_i \in \Delta_i$ . As a result, the new population profile for population  $i$  becomes  $x_i^\epsilon = \epsilon \cdot y_i + (1 - \epsilon)x_i$ . If we denote the set of payoff functions for all populations as  $u = \{u_1, \dots, u_n\}$  where the payoff of an agent in population  $i$  taking strategy  $s_j^i$  is  $u_i(s_j^i, x_{-i})$ . The average payoff of the rational players in the population  $i$  before and after the invasion can be expressed as  $u_i(x_i, x_{-i}) = \sum_{j=1}^{m_i} x_i(j)u_i(s_j^i, x_{-i})$  and  $u_i(x_i, x_{-i}^\epsilon) = \sum_{j=1}^{m_i} x_i(j)u_i(s_j^i, x_{-i}^\epsilon)$  respectively and the average payoff of the mutant invaders in population  $i$  is  $u_i(y_i, x_{-i}^\epsilon) = \sum_{j=1}^{m_i} y_i(j)u_i(s_j^i, x_{-i}^\epsilon)$ . Now, we can describe a general mechanism as  $\mathcal{M} = \{I, \Theta, u\}$  and the robustness of a desired equilibrium can be defined as follows.

**Definition 1** (Robustness of a desired equilibrium) Given an incentive mechanism  $\mathcal{M} = \{I, \Theta, u\}$  with a desired equilibrium  $x \in \Theta$ , the robustness of  $x$  is  $R$  such that

$$R = \min_{y \in \Theta(y \neq x)} \arg \max_{\epsilon \in (0,1)} u_i(x_i, x_{-i}^\epsilon) > u_i(y_i, x_{-i}^\epsilon), \quad \forall i \in I. \tag{1}$$

In other words, the robustness of a desired equilibrium is the maximum proportion of invaders such that the desired equilibrium strategy is still the best strategy for rational players in each population.

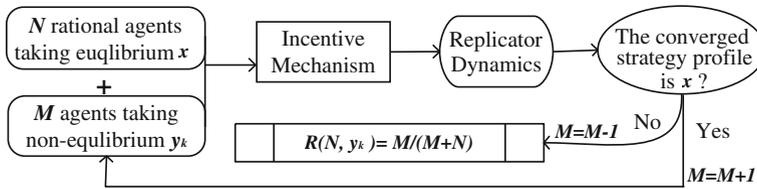
**Observation 1** *The robustness of a strictly dominant equilibrium is 1.*

*Proof* Given that the desired equilibrium  $x$  is a strictly dominant equilibrium,  $u_i(x_i, x_{-i}^\epsilon) > u_i(y_i, x_{-i}^\epsilon)$  in Eq.(1) would be always true for all  $\epsilon \in (0, 1)$ . Thus the robustness of the equilibrium is 1, i.e.,  $R = 1$ .

However, the computational cost of formally analyzing the robustness measure in Definition 1 makes such analysis infeasible. First, the size of the population corresponding to each player position is assumed to be infinite in the definition so as to model all possible mixed strategies of the player. The payoff of a population is thus difficult to calculate through aggregating all individuals’ utilities. Second, the strategy of the mutants  $y \neq x$  is continuous and it may be impossible to verify Eq.(1) for all  $y \in \Theta$ . Moreover, the size of players  $I$  in a realistic game can be very large, making the robustness measure computationally prohibitive. Given the difficulty of formally analyzing the robustness of incentive mechanisms, we consider finite populations to be able to realistically measure the robustness of incentive mechanisms.

Moreover, the mutants or invaders in a realistic mechanism always attack the mechanism using several mature techniques or strategies, e.g., constant attack and whitewashing attack in reputation systems [12]. It is also observed that these attacks are typically conducted by a single player. Thus, it is practical to measure the robustness of an incentive mechanism against representative types (and their combinations) of mutants in a single population.

<sup>2</sup> When the game is symmetric and the positions of the  $n$  players are equivalent, a single population could represent all the players who interact with each other.



**Fig. 1** Evaluating robustness based on replicator dynamics

We note that the robustness measure in Definition 1 is a static concept. A mechanism is always implemented in real time and players could dynamically adjust their strategies as time goes on, which has been modeled by the well-known replicator dynamics (RD) model [4, 8, 17, 18]—a representative model for evolutionary dynamics (population dynamics)—for capturing long-range evolutionary dynamics with natural selection forces. If agents using the desired equilibrium strategy could achieve higher payoff (fitness) in a population, then they are more likely to reproduce and produce offsprings which use the same strategy (replicator) and eventually through this process the strategy distribution of the population will converge to a stable equilibrium profile. Based on this analysis, we now propose an evolutionary dynamics based practical robustness measure that evaluates the robustness of an equilibrium against any specific non-equilibrium strategy.

In Fig. 1, agents adopting the non-equilibrium strategy  $y_k \neq x_k$  are gradually added to the original population of rational agents (population without mutants) and the resultant populations repeatedly interact under the mechanism. The strategy distribution of the rational agent populations evolves through the RD model until convergence. Successively more mutants are added until the population converges to a distribution different from the desired equilibrium.

**Definition 2** (Robustness against a non-equilibrium strategy) The robustness of an incentive mechanism  $\mathcal{M}$ , with desired equilibrium strategy profile  $x$ , against a non-equilibrium strategy  $y_k$ , used by bounded-rational invaders, is

$$R(N, y_k) = \frac{M}{N + M} \tag{2}$$

where  $N$  is the finite population size for each player position and  $M$  is the maximum number of bounded rational invaders such that the converged strategy profile, under the RD model, does not deviate from the desired equilibrium.

We assume that each player position has the same population size and the robustness of an incentive mechanism is a function of the population size and the non-equilibrium strategy. We will evaluate the robustness of a mechanism with different population sizes and against different non-equilibrium strategies.

### 4 Simulation framework

In this section, we propose a simulation-based framework to calculate the robustness of an incentive mechanism. Our framework is based on replicator dynamics [16] and is used to study the dynamics of the equilibrium strategy distribution in populations that include agents using a non-equilibrium strategy. Strategies which can garner higher payoffs become more prevalent while ineffective strategies are gradually eliminated through the replicator dynamics

process. According to Definition 2, we gradually add agents using a specific non-equilibrium strategy and let the populations evolve over time until convergence. We stop inserting more agents when the equilibrium strategies are abandoned by the populations of rational agents in the converged population profile after a sufficiently large number of evolutionary dynamics steps.

Recall that an incentive mechanism is denoted by  $\mathcal{M}(I, \Theta, u)$ . The population size for every player position is  $N$ , and the pure strategy set of player  $i \in I$  is  $S_i = \{s_1^i, \dots, s_{m_i}^i\}$ . The desired equilibrium is  $x = (x_1, \dots, x_n)$  where  $x_i$  is the desired equilibrium strategy of population  $i$ . The robustness of a mechanism  $\mathcal{M}$  against a non-equilibrium strategy  $y_k$  is evaluated using Algorithm 1.

```

Input :  $\mathcal{M}$  : incentive mechanism under evaluation;
          $n$  : number of players in  $\mathcal{M}$ ;
          $N$  : population size for each player;
          $x$  : desired equilibrium;
          $T$  : number of evolutionary process repeated;
          $y_k$  : non-equilibrium strategy of invading agents ;
Output:  $R(N, y_k)$ , Robustness of  $\mathcal{M}$  against  $y_k$ ;
1 Initialize  $\varepsilon = \frac{1}{T}$ ;  $\xi = 0$ ;
2 Set  $M = 1$ ;
3 repeat
4   For each player  $j$ ,  $N$  agents taking strategy  $x_j$ ;
5   Add  $M$  agents taking strategy  $y_k$  to population  $k$ ;
6    $\lambda = 0$ ;
7   for  $Times = 1 \rightarrow T$  do
8      $\bar{x} = x$ ;
9     while  $\bar{x}$  is not converged do
10      Run  $\mathcal{M}$  in the  $n$  populations;
11      for  $i = 1 \rightarrow n$  do
12        Calculate average payoff  $u_i(s_j^i, \bar{x}_{-j})$  for each strategy  $s_j^i \in S_i$ ;
13        Generate new population  $i$ ;
14       $\bar{x}$  = the strategy profile of rational populations ;
15      if  $\bar{x} \neq x$  then
16         $\lambda = \lambda + 1$ ;
17     $\xi = \frac{\lambda}{T}$ ;  $M = M + 1$ ;
18 until  $\xi \geq \varepsilon$ ;
19 return  $R(N, y_k) = \frac{M-1}{N+M-1}$ ;

```

**Algorithm 1:** Evaluating the Robustness of an Incentive Mechanism  $\mathcal{M}$  Against a Non-Equilibrium Strategy  $y_k$ .

In Algorithm 1, we first set three parameters which are used in the simulation framework (Line 1). We initialize the number of agents taking non-equilibrium strategy to be  $M = 1$  (Line 2), then gradually increase the number  $M$  until the rational agents deviate from their desired equilibrium profile  $x$  in at least  $\varepsilon = \frac{1}{T}$  (Line 4–18). Due to the randomness involved in the evaluation framework (Line 11–14), we repeat the evolutionary process for  $T$  times and  $T = 50$  when we implement it to evaluate the robustness of incentive mechanisms in e-marketplaces in Sect. 6, and in each round, the strategy profile of each population  $i$  evolves (Line 9–14) until it converges. The “ $\bar{x}$  is not converged” (Line 9) in the framework means that

the strategy profile  $\bar{x}$  is still not stable enough according to the convergence criteria specified in Sect. 5.2. The variable  $\lambda$  records the number of times out of  $T$  rounds that the rational agents deviate from their equilibrium strategy ( $\bar{x} \not\approx x$ ), i.e., the converged strategy profile deviates from the desired equilibrium  $x$  as indicated by  $\max_{i \in I, j \in \{1, \dots, m_i\}} |\bar{x}_i(j) - x_i(j)| > 0.01$  (Line 15–16). After completing  $T$  repetitions, we calculate the percentage of times that the populations of rational agents deviate from the desired equilibrium (Line 17). We use the probability  $\varepsilon = \frac{1}{T}$  to specify the extent to which the deviation is considered acceptable, which is further explained in Proposition 2. If  $\xi < \varepsilon$ , we conclude that the populations of rational agents could still sustain their desired equilibrium when invaded by  $M$  agents using the non-equilibrium strategy  $y_k$ , and we will then repeat the number of simulation runs after increasing  $M$  by 1. When  $y_k$  is a mixed strategy, the  $M$  agents probabilistically take each pure strategy according to the probability distribution stated in  $y_k$ . This continues until the deviation of the rational agent strategies in the converged population becomes too frequent ( $\xi \geq \varepsilon$ ). Then, the robustness of the mechanism against  $y_k$  is returned as  $R(N, y_k) = \frac{M-1}{N+M-1}$  (Line 19), as per Definition 2.

The new population for each player (Line 11–14) will be generated after each evolutionary dynamics step based on the Moran process [16, 17], a well-known evolutionary dynamics process for finite populations. More specifically, we divide each population  $i$  with  $N$  rational agents into  $k_i$  subsets where agents in each set  $j$  take pure strategy  $s_j^i \in S_i$ . The size of the corresponding subsets are  $N_1, \dots, N_{m_i}$ , and the average payoff of the set  $j$  is  $u_i(s_j^i, \bar{x}_{-i})$ . If subset  $j$  is empty, i.e.,  $N_j = 0$ ,  $u_i(s_j^i, \bar{x}_{-i}) = 0$ . In each step, we randomly choose an agent from each population  $i$  to be replaced by a new agent using strategy  $s_j^i$  with probability

$$P_i(s_j^i) = \frac{N_j \times u_i(s_j^i, x_{-i})}{\sum_{l=1}^{m_i} N_l \times u_i(s_l^i, x_{-i})} \tag{3}$$

and the other  $N - 1$  agents in the new generation inherit the strategies of the  $N - 1$  agents (excluding the selected agent) one by one. According to Eq. (3), a strategy  $s_j^i$  in population  $i$  is more likely to be taken by more individuals in the next generation when the proportion of agents taking  $s_j^i$  is larger and average payoff  $u_i(s_j^i, x_{-i})$  is higher.

Note that a strategy  $s_j^i$  that becomes extinct ( $N_j$  drops to 0) will never re-appear as the probability of another agent adopting that strategy  $P_i(s_j^i)$  also becomes 0. To avoid the system being locked in initial homogeneous states with all agents using the same desired strategy, it is necessary to use variation to re-introduce extinct strategies.

For each population  $i$ , if the strategy set  $S_i^d = \{s_j^i | N_j = 0, j = 1, \dots, m_i\}$  is not empty, the framework will run the variation process. First, we randomly choose one strategy  $s_j^i$  appearing in the population with probability  $P_i(s_j^i)$ , and uniformly choose another strategy  $s_{j'}^i$  in the extinct strategy set  $S_i^d$ . We then change the strategy of an agent from subset  $j$  to be  $s_{j'}^i$  with probability

$$F_i(s_j^i \rightarrow s_{j'}^i) = \frac{1}{e^{\rho(u_i(s_j^i, x_{-i}) - u_i(s_{j'}^i, x_{-i}))} + 1} \tag{4}$$

where  $\rho$  is a positive constant, representing the intensity of variation [17, 26]. When  $\rho \ll 1$ ,  $F_i(s_j^i \rightarrow s_{j'}^i)$  converges to 0.5 which produces high neutral fluctuation. When  $\rho \rightarrow +\infty$ ,  $F_i(s_j^i \rightarrow s_{j'}^i)$  becomes a step function where the probability of the strategy  $s_{j'}^i$  appearing is 1 if  $u_i(s_j^i, x_{-i}) < u_i(s_{j'}^i, x_{-i})$ , and 0 otherwise. We have conducted extensive

evaluations with different  $\rho$  values: the results show that  $\rho$  does not significantly influence the evolution of the population profile. We have verified that the evaluated robustness is not sensitive to the value of  $\rho$  in the Prisoner’s Dilemma game, as presented in Fig. 4. This is due to the fact that a strategy could only reappear in the population through the variation process when the strategy has already been extinct. Therefore, we use  $\rho = 1$  in all our simulations unless otherwise specified.

The robustness of mechanisms is related to the population size  $N$ . However, for special settings, the robustness  $R(N, y_k)$  can be independent of population size  $N$  as shown in Proposition 1. Generally speaking,  $R(N, y_k)$  increases as the population size grows because the stochastic effect is less impactful in a larger population which results in a larger robustness value. The value  $R(N, y_k)$  finally converges when the population size is large enough, which will be evaluated in Proposition 2.

**Proposition 1** *There exists  $\varepsilon = G(N)$  where  $\frac{dG(N)}{dN} < 0$  such that the evaluated robustness of a mechanism is independent of the population size  $N$ .*

*Proof* We first prove this result in a special case and then show how to generalize it. Consider a special mechanism for a symmetric game in which each player has two pure strategies  $s_1$  and  $s_2$ . The desired equilibrium strategy is pure strategy  $s_1$  denoted by  $x = [1, 0]$ . The other strategy  $s_2$  is  $y = [0, 1]$ . The population containing  $N$  rational agents whose strategy profile is  $x$  and  $M$  bounded rational agents whose strategy profile is  $y$  such that the expected utility of agents taking  $s_1$  in the rational population is  $a = u(x, \delta \cdot y + (1 - \delta)x)$  and the utility of taking  $s_2$  is  $b = u(y, \delta \cdot y + (1 - \delta)x)$  for any  $\delta \in (0, 1)$ .

Among the rational agents, the number of agents taking  $s_1$  is denoted by  $N_1$ , and the number of agents taking  $s_2$  by  $N_2$ . Initially,  $N_1 = N$  and  $N_2 = 0$ . According to the Moran process in generating a new population, only one agent could change its strategy in each step. Suppose that the probability of the number of desired strategy taker  $N_1$  changing from  $i$  to  $j$  after each step is  $T_{i,j}$ .  $T = \{T_{i,j}\}$  is a  $(N + 1) \times (N + 1)$  tridiagonal matrix where

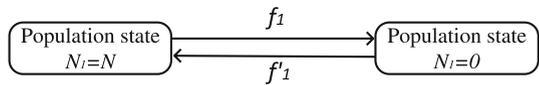
$$\begin{aligned} T_{0,1} &= 0 \\ T_{0,0} &= 1 \\ &\dots \dots \\ T_{i,i-1} &= \frac{i}{N} \cdot \frac{(N-i)b}{i \cdot a + (N-i)b} \\ T_{i,i+1} &= \frac{N-i}{N} \cdot \frac{i \cdot a}{i \cdot a + (N-i)b} \\ T_{i,i} &= 1 - T_{i,i-1} - T_{i,i+1} \\ &\dots \dots \\ T_{N,N-1} &= 0 \\ T_{N,N} &= 1 \end{aligned}$$

Next we show that when  $\varepsilon = G(N) = \frac{1}{N+1}$ , the robustness of the mechanism is independent of the population size  $N$ . The dynamic process has two absorbing states (when the population profile becomes one of these, the Moran process will converge):  $N_1 = 0$  and  $N_1 = N$ . We now calculate the probability,  $f_i$ , of the population converging to state  $N_1 = 0$  when starting at state  $N_1 = i$ , which is recursively defined as,  $\forall i \in \{1, \dots, N - 1\}$ ,

$$f_i = T_{i-1,i} \cdot f_{i-1} + T_{i,i} \cdot f_i + T_{i+1,i} \cdot f_{i+1}. \tag{5}$$

We note that  $f_0 = 1$  and  $f_N = 0$ , and hence  $f_1, \dots, f_{N-1}$  can be calculated by solving these  $N - 1$  linear equations [17]:

**Fig. 2** The state transition of the population



$$f_i = \begin{cases} \frac{i}{N} & \eta = 1 \\ \frac{1 - (\frac{1}{\eta})^i}{1 - (\frac{1}{\eta})^N} & \eta > 1 \end{cases} \tag{6}$$

where  $\eta = \frac{b}{a}$ . When  $\eta \rightarrow 1$ , we can derive that  $\lim_{\eta \rightarrow 1} f_i = \frac{i}{N}$ , and it is how we obtain the value  $f_i$  for  $\eta = 1$ . Since  $f_1$  and  $f_{N-1}$  will be utilized in the following proof, we specify them as follows:

$$f_1 = \begin{cases} \frac{1}{N} & \eta = 1 \\ \frac{1 - \frac{1}{\eta}}{1 - (\frac{1}{\eta})^N} & \eta > 1 \end{cases} \quad \text{and} \quad f_{N-1} = \begin{cases} \frac{N-1}{N} & \eta = 1 \\ \frac{1 - (\frac{1}{\eta})^{N-1}}{1 - (\frac{1}{\eta})^N} & \eta > 1 \end{cases} \tag{7}$$

The value  $f_1$  is also called the *fixation probability* that a single, bounded rational agent using  $y$  makes all the other  $N - 1$  rational agents using  $x$  to eventually choose strategy  $y$ . Conversely, the probability that a single agent using  $x$  makes all the other  $N - 1$  agents using  $y$  to eventually use strategy  $x$ , is denoted by  $f'_1$  which is calculated by  $1 - f_{N-1}$ .

We could derive that  $f_1 < \frac{1}{N}$  when  $\eta < 1$ , and  $\frac{\partial f_1}{\partial \eta} > 0$ .<sup>3</sup> Since  $\lim_{\eta \rightarrow 1} \frac{1 - \frac{1}{\eta}}{1 - (\frac{1}{\eta})^N} = \frac{1}{N}$  when  $N$  is a constant, and  $\frac{\partial[(1 - \frac{1}{\eta}) / (1 - (\frac{1}{\eta})^N)]}{\partial \eta} > 0$ ,  $f_1$  is a strictly increasing function of  $\eta$ . Similarly, we could derive that  $f'_1$  is a strictly decreasing function of  $\eta$ . Intuitively, as  $\eta = \frac{b}{a}$  increases, the strategy  $y$  can obtain higher utility compared to strategy  $x$ , making it more likely that all agents take  $y$ , i.e.,  $f_1$  increases with  $\eta$ , and correspondingly, it is less likely that all agents adopt  $x$ , i.e.,  $f'_1$  decreases with  $\eta$ .

Considering the variation process, we assume a strategy  $x$  (or  $y$ ) can be re-introduced into the population even after  $x$  (or  $y$ ) became extinct. For a population where all agents use  $x$  at state  $N_1 = N$ , the probability of it evolving to the state  $N_1 = 0$  (considering one mutator) is  $f_1$  and the probability of the population evolving from  $N_1 = 0$  to  $N_1 = N$  is  $f'_1$ . This state transition is described as in Fig. 2. Then, the probability that the population where all agents use  $x$  evolve to a population where all agents use  $y$  is

$$\begin{aligned} \xi(N, \eta) &= f_1 - f_1 f'_1 + f_1(f_1 f'_1) - f_1 f'_1(f_1 f'_1) + \dots \\ &= f_1 + f_1(f_1 f'_1) + f_1(f_1 f'_1)^2 \dots \\ &\quad - f_1 f'_1 - f_1 f'_1(f_1 f'_1) - f_1 f'_1(f_1 f'_1)^2 - \dots \\ &= \frac{f_1 - f_1 f'_1}{1 - f_1 f'_1} \end{aligned} \tag{8}$$

From Eqs. (7) and (8) and using  $\lim_{\eta \rightarrow 1} f'_1 = \frac{1}{N}$ , we have  $\xi(N, \eta)|_{\eta=1} = \frac{1}{N+1}$ . We can also derive that

$$\frac{\partial \xi(N, \eta)}{\partial \eta} = \frac{\frac{\partial f_1}{\partial \eta}(1 - f'_1) - \frac{\partial f'_1}{\partial \eta} f_1(1 - f_1)}{(1 - f_1 f'_1)^2} > 0 \tag{9}$$

<sup>3</sup> We do not consider the case of  $\eta < 1$  in Eq.(7) as in that case  $f_1 < \frac{1}{N}$ , which is a very small value for a large population size and hence equilibrium strategies are very likely to be robust against one or few invaders.

In Algorithm 1,  $\varepsilon$  represents the maximum acceptable rate of deviation from the desired strategy. The algorithm increases  $M$  until  $\xi(N, \eta) = \varepsilon$ . If  $\varepsilon = \frac{1}{N+1}$ , the condition becomes  $\xi(N, \eta) = \frac{1}{N+1}$ . Since  $\xi(N, \eta)$  is a strictly increasing function of  $\eta$ , it is satisfied only when  $\eta = 1$ , i.e.,  $b = a$ , which is independent of the population size  $N$ . Therefore, in the above special case, there exists a function  $G(N) = \frac{1}{N+1}$  such that the robustness is independent of the population size  $N$ , as  $\varepsilon = G(N)$  where  $\frac{dG(N)}{dN} < 0$ .

When the population have more than two pure strategies, the robustness of the mechanism will still be independent of  $N$  by treating  $s_2$  in the above proof as the non-equilibrium pure strategy such that the agents taking the non-equilibrium strategy could achieve the highest payoff among all the non-equilibrium pure strategies, and  $\eta = \frac{b}{a}$  where  $b$  is the payoff of taking this strategy and  $a$  is the payoff of taking the desired equilibrium strategy. For a general mechanism with  $n > 2$  players, it can be easily derived that such  $G(N)$  still exists, since the population for each player could evolve in the same way as that in the special case.

However, it is more reasonable to set  $\varepsilon$  as a constant since  $\varepsilon$  represents the acceptable rate of deviation from the desired strategy.

**Proposition 2** *When  $\varepsilon = \varepsilon_0$  is a constant, the robustness of a mechanism increases with the population size  $N$  where  $N > \frac{1}{\varepsilon_0} + 1$ , and finally converges at a value when the agents of a population taking a non-equilibrium strategy can achieve  $\frac{\varepsilon_0}{1-\varepsilon_0}$  times payoff more than that of the desired equilibrium strategy.*

*Proof* In Algorithm 1, we stop increasing  $M$  when the rate of deviation from the desired strategy,  $\xi(N, \eta) = \frac{\lambda}{T}$ , becomes more than  $\varepsilon$ . For a constant  $\varepsilon_0$ , if the robustness of a mechanism increases with the population size, then  $\xi(N, \eta)$  should be a decreasing function of the population size  $N$ .

Given that  $N > \frac{1}{\varepsilon_0}$ , we can obtain  $\eta > 1$  when  $\xi(N, \eta) = \varepsilon_0$  since  $\frac{\partial \xi(N, \eta)}{\partial \eta} > 0$  and  $\xi(N, 1) = \frac{1}{N+1} < \varepsilon_0$ . Since we have (according to Eq. 6)

$$\begin{aligned} \frac{\partial f_1}{\partial N} &= \frac{(\frac{1}{\eta})^N \ln(\frac{1}{\eta})(1 - \frac{1}{\eta})}{(1 - (\frac{1}{\eta})^N)^2} = \frac{\eta^{N-1} \ln(\eta)(1 - \eta)}{(1 - \eta^N)^2} \\ \frac{\partial f'_1}{\partial N} &= \frac{\eta^N \ln(\eta)(1 - \eta)}{(1 - \eta^N)^2} \end{aligned} \tag{10}$$

then  $\frac{\partial f'_1}{\partial N} = \eta \frac{\partial f_1}{\partial N} < 0$ . According to Eq. (8), we derive that

$$\begin{aligned} \frac{\partial \xi(N, \eta)}{\partial N} &= \frac{\frac{\partial f_1}{\partial N}(1 - f'_1) - \frac{\partial f'_1}{\partial N} f_1(1 - f_1)}{(1 - f_1 f'_1)^2} \\ &= \frac{1}{(1 - f_1 f'_1)^2} \frac{\partial f_1}{\partial N} (1 - f'_1 - \eta f_1 + \eta f_1^2) \\ &= \frac{(1 - f_1)(1 - f'_1)}{(1 - f_1 f'_1)^2} \frac{\partial f_1}{\partial N} < 0. \end{aligned} \tag{11}$$

Thus, for a constant  $\varepsilon_0$ , the robustness is a decreasing function of the population size  $N$  where  $N > \frac{1}{\varepsilon_0} + 1 > 1$ .

We next show the convergence of the robustness value as the population size approaches infinity. Since  $\frac{\partial \xi(N, \eta)}{\partial N} < 0$  and  $\lim_{N \rightarrow +\infty} \frac{\partial \xi(N, \eta)}{\partial N} = 0$ , then  $\xi(N, \eta)$  is a decreasing function of

$N$  and  $\lim_{N \rightarrow +\infty} \xi(N, \eta) = 1 - \frac{1}{\eta}$ . Given that  $\varepsilon = \varepsilon_0$ , we could derive that  $\eta = \frac{1}{1-\varepsilon_0}$ . In other words, the robustness of an incentive mechanism will converge to a constant as population size approaches infinity when  $\eta$  satisfies the condition  $\eta = \frac{1}{1-\varepsilon_0}$ . This means that there exists a non-equilibrium strategy whose payoff is  $\frac{\varepsilon_0}{1-\varepsilon_0}$  times more than that of the desired equilibrium strategy.

We have shown that the robustness of an incentive mechanism is an increasing function of the population size and is bounded from below, provided that  $\varepsilon$  is a constant. As  $\varepsilon$  is set to a very large value (e.g., approaching to 1 where most agents take the non-equilibrium strategy), then the framework would allow the non-equilibrium to achieve a payoff of  $\frac{\varepsilon_0}{1-\varepsilon_0} \rightarrow +\infty$  times that of the equilibrium strategy when the population is large. To avoid this singularity in the following simulations, we set  $\varepsilon = \frac{1}{T} = 0.02$  which is a sufficiently small but measurable value<sup>4</sup> and will gradually increase the population size until the robustness value converges. We observe that by setting  $\varepsilon$  to a small value the evaluation framework will converge at nearly the same robustness value as the ones analyzed by the extended simplex analysis in the next section.

### 5 Validation of the framework on representative games

In this section, we validate the proposed simulation framework by comparing the simulated results with the analytical ones in simple representative games. The simplex analysis approach is used to analytically evaluate the robustness of equilibrium strategies.

#### 5.1 Extended simplex analysis approach

For ease of presentation, we consider a two-player asymmetric game<sup>5</sup> to introduce the simplex analysis which is a well-known approach for analyzing evolutionary dynamics [20], and we extend it to calculate the robustness of some representative games.

The strategy sets of the two players are  $S_1 = \{s_1^1, \dots, s_{m_1}^1\}$  and  $S_2 = \{s_1^2, \dots, s_{m_2}^2\}$ , and their payoff matrices are  $A_1$  and  $A_2$  which are  $m_1 \times m_2$  and  $m_2 \times m_1$  matrix, respectively:

$$A_1 = \begin{bmatrix} a_{11}^1 & \cdots & a_{1m_2}^1 \\ \vdots & \ddots & \vdots \\ a_{m_1 1}^1 & \cdots & a_{m_1 m_2}^1 \end{bmatrix} \quad A_2 = \begin{bmatrix} a_{11}^2 & \cdots & a_{1m_1}^2 \\ \vdots & \ddots & \vdots \\ a_{m_2 1}^2 & \cdots & a_{m_2 m_1}^2 \end{bmatrix} \tag{12}$$

The population profile of population 1 is  $x_1 = [x_1(1), \dots, x_1(m_1)]$  where  $x_1(j)$  is the proportion of agents in this population using strategy  $s_j^1 \in S_1$  and  $\sum_{j=1}^{m_1} x_1(j) = 1$ . Similarly, the population profile of population 2 is  $x_2 = [x_2(1), \dots, x_2(m_2)]$ .

The evolutionary dynamics of the populations is then the dynamics of  $x_1$  and  $x_2$ . Since the population profile  $x_i$  for population  $i$  is a strategy distribution satisfying  $\sum_{j=1}^{m_i} x_i(j) = 1$ ,  $x_i$  is also called as a *simplex*. This analysis approach is thus called the *simplex* analysis. The evolutionary dynamics of the simplex is described by the replicator dynamics [17,28]:

<sup>4</sup> There are two considerations for choosing an appropriate value for the parameter  $\varepsilon$ : (1)  $\varepsilon$  should be no less than its measurable accuracy  $\frac{1}{T}$ ; (2) it is better to use a smaller value for  $\varepsilon$ .

<sup>5</sup> For a symmetric game played in a single population, it can be considered as two replicated populations.

$$\nabla x_1 = x_1 \left[ \mathbf{A}_1 x_2 - \frac{x_1^T \mathbf{A}_1 x_2}{m_1} \right] \tag{13}$$

$$\nabla x_2 = x_2 \left[ \mathbf{A}_2 x_1 - \frac{x_2^T \mathbf{A}_2 x_1}{m_2} \right] \tag{14}$$

where  $\nabla x_1$  and  $\nabla x_2$  denote the gradient vector of the population profiles at  $x_1$  and  $x_2$ , respectively. These equations show that the proportion of agents taking a certain pure strategy in a population increases if and only if the agents could achieve higher utility than the average payoff of the population, and vice versa.

According to Eqs. (13) and (14), we could calculate the gradient at each point in the simplex space. Starting from a particular point, we generate a smooth trajectory by moving a small distance along the calculated gradient, until reaching a fixed point which is called an *attractor*. Each attractor attracts the points around it.

Let  $x_i^*$  be an attractor of the population  $i$  where  $i = 1$  or  $2$ . Since the attractor is the point to which the points around it will converge, the following equations should be satisfied:

$$\begin{aligned} \nabla x_i &= \mathbf{0}_{1 \times m_i}, & x_i &= x_i^* \\ \nabla x_i \cdot (x_i^* - x_i) &> 0, & x_i &\in x_i^* \pm \epsilon \end{aligned} \tag{15}$$

where  $x_i^* \pm \epsilon$  represents the point set around  $x_i^*$  in the space.

Next, we extend this approach to analytically evaluate the robustness of an incentive mechanism. One population is denoted by  $i$  and the invaders are added to population  $i$  taking non-equilibrium strategy  $y_i$ , and the other population is denoted by  $j$ . We evaluate the maximum proportion of invaders such that the population profile of rational agents  $x$  is an attractor  $x^*$ . Based on this intuition, we make the following modifications to the basic simplex analysis approach:

1. Update the population profile of population  $i$  from  $x_i$  to

$$\hat{x}_i = \{(1 - R)x_i, R \cdot y_i\} \tag{16}$$

by separately listing the strategy profiles of rational agent population  $i$  and the bounded rational strategy, where  $R$  is the proportion of the bounded rational agents. The desired equilibrium<sup>6</sup> becomes  $\hat{x}_i^* = \{(1 - R)x_i^*, R \cdot y_i\}$ .

2. Update the payoff matrix of the population  $j$  ( $j \neq i$ ) from  $\mathbf{A}_j$  to  $\hat{\mathbf{A}}_j$  by adding  $m_i$  columns to describe the payoff of the population  $j$  interacting with the bounded rational agents taking  $y_i$ . Then  $\hat{\mathbf{A}}_j$  is a  $m_j \times 2m_i$  matrix.
3. Calculate  $\nabla \hat{x}_i = [\nabla x_i, \mathbf{0}_{1 \times m_1}]$ , and  $\nabla x_j = x_j [(\hat{\mathbf{A}}_j \hat{x}_i - \frac{x_j^T \hat{\mathbf{A}}_j \hat{x}_i}{m_j})]$ .
4. Calculate solution set  $\Omega = \{R\}$  of Eq. (16) satisfying:

$$\nabla \hat{x}_i = \mathbf{0}_{1 \times 2m_i}, \quad \nabla x_j = \mathbf{0}_{1 \times m_j} \tag{17}$$

and one of the following  $m_i + m_j$  equations

$$\begin{aligned} \frac{\partial \nabla \hat{x}_i}{\partial \hat{x}_i^{(k)}} &= 0, & \text{where } k &= 1, \dots, m_i \\ \frac{\partial \nabla x_j}{\partial x_j^{(k)}} &= 0, & \text{where } k &= 1, \dots, m_j \end{aligned} \tag{18}$$

5. Calculate the robustness value:

$$R(\mathcal{M}, y_i) = \begin{cases} \max\{\Omega\} & \Omega \neq \emptyset \\ 1 & \Omega = \emptyset \end{cases} \tag{19}$$

<sup>6</sup> We assume that the strategy of bounded rational agents cannot evolve as they are not rational.

The robustness computed using the simplex analysis is exactly the one defined in Definition 2 when the population size is infinite. Thus, our simulated robustness should converge to the analytical value as population size  $N$  becomes sufficiently large.

It worth noticing that the extended simplex analysis approach can only be used to evaluate the robustness of stylized games which have explicit payoff matrices. For complex real mechanisms, the payoff matrix is not explicitly available or even cannot be estimated through empirical approaches, such as [31,33]. As a result, this analytical approach may not be generally applicable. Thus, we validate the proposed simulation frame on several representative games which can be analyzed using the extended simplex analysis approach.

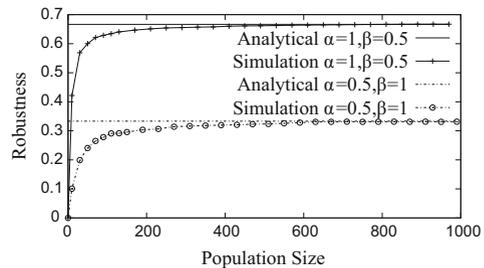
### 5.2 Evaluation on representative games

In this section, we evaluate the proposed simulation framework by comparing the simulated results with the analytical ones calculated by the simplex analysis approach in some representative games.

We first discuss the parameters used in the simulation framework. We conduct the experiments with the *Symmetric Coordination Game* with population sizes  $N = 400, 600,$  and  $800,$  respectively. In the game, the pure strategy set for each player is  $\{Left, Right\}$  and the payoff matrix is shown at the left side of Fig. 3. Both players gain the same payoff if they choose the same strategy, i.e.,  $\alpha$  if both choose *Left* and  $\beta$  if both choose *Right*. Thus, two pure strategy Nash equilibria exist. Assume that  $\{Left, Left\}$  is the desired equilibrium, rational agents would always choose *Left* given that other agents are also rational and choose *Left*, i.e.,  $x^* = \{1, 0\}$ . The bounded rational agents, instead, would choose any other strategy but *Left*. We let the bounded rational agents always choose *Right*, i.e.,  $y = \{0, 1\}$ , and both  $\alpha$  and  $\beta$  be 1, in this experiment. The evaluated robustness values are listed in Table 1.

For the same population size (e.g.,  $N = 400$ ), the robustness of the mechanism increases as  $\epsilon$  increases, and for the same  $\epsilon$  the robustness value increases as the population size increases as stated in Proposition 2. In the following evaluation framework, we set  $\epsilon = \frac{1}{T} = 0.02$ . Moreover, to justify that a population profile has converged, we set the convergence criterion

		Player II	
		Left	Right
Player I	Left	$\alpha$	0
	Right	0	$\beta$



**Fig. 3** The robustness of the Symmetric Coordination Game

**Table 1** The robustness of the coordinated game ( $\alpha = 1, \beta = 1$ ) with different  $N$  and  $\epsilon$

$N$	$\epsilon = 0.1$	$\epsilon = 0.3$	$\epsilon = 0.9$
400	$325/(325 + N) = 0.448$	$330/(330 + N) = 0.452$	$375/(375 + N) = 0.484$
600	$525/(525 + N) = 0.467$	$532/(532 + N) = 0.470$	$576/(576 + N) = 0.490$
800	$725/(725 + N) = 0.475$	$735/(735 + N) = 0.479$	$779/(779 + N) = 0.493$

		Player II	
		Cooperation	Defection
Player I	Cooperation	b	d
	Defection	d	c

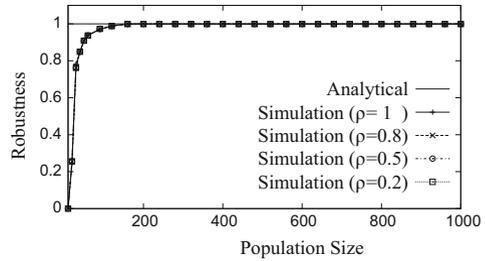


Fig. 4 The robustness of the Prisoner’s Dilemma

		Woman	
		Boxing	Ballet
Man	Boxing	$\alpha$	$\beta$
	Ballet	$0$	$\alpha$

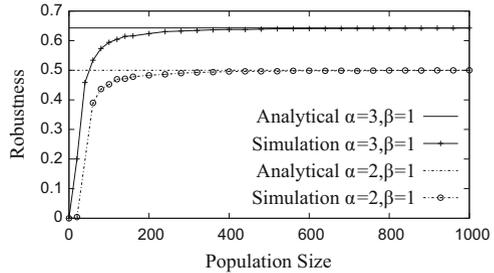


Fig. 5 The robustness of the Battle of Sexes game

$\max_{j=1}^L (\max |\bar{x}(t) - \bar{x}(t-j)|) \leq \vartheta$  where  $t > L$  and  $\bar{x}(t)$  is the population profile of the population at time step  $t$ , which includes two parameters  $L$  and  $\vartheta$ . This criterion specifies that the maximum difference of the converged profile of the population at step  $t$  with that over the previous  $L$  time steps is no more than a threshold  $\vartheta$ . We use  $L = 100$  and  $\vartheta = e^{-2}$  in our experiments.

The first game we evaluated is a *Symmetric Coordination Game* as described earlier.

According to Eqs. (17)–(19), The robustness of the game against this non-equilibrium strategy is  $\frac{\alpha}{\alpha+\beta}$ . We implement this game in our simulation framework by varying  $\alpha$  and  $\beta$  values. Figure 3 shows that the simulated robustness does indeed always converge to the analytical results as population size is sufficiently large.

The second game we consider is the *Prisoner’s Dilemma*, where each player has two pure strategies: *Cooperation* and *Defection*. Its payoff matrix is shown in the matrix on the left-hand side of Fig. 4, where  $d > b > c > a$ ,  $2b > a + d$ .

The population profile is  $x = (x(1), x(2))$  where  $x(1)$  denotes the proportion of population taking strategy *cooperation*. In this game, there is one unique dominant Nash equilibrium where all players take pure strategy *defection*, denoted by  $x^* = (0, 1)$ . The non-equilibrium strategy is *cooperation*. Based on Observation 1 and Eqs. (17) and (18), the robustness of the game against this non-equilibrium strategy is 1. Figure 4 shows that the simulated robustness converges to 1. We also set  $\rho$  with different values (0.2, 0.5, 0.8 and 1), and the simulation results are exactly the same (the curves for different  $\rho$  overlap with each other), showing that our simulation results are not sensitive to the parameter  $\rho$ .

The third game used for validating our approach is the *Battle of Sexes game*. There are two players: *Man* and *Woman*, and a man or woman has two choices (strategies): *Boxing* and *Ballet*. The payoff matrix of the players is shown in the matrix on the left side of Fig. 5.

A man prefers to watch *Boxing* and a woman prefers *Ballet*. In this game, a pair of individuals is randomly chosen to show up at the sites for either boxing or ballet. If the

two individuals have the same sex or different sites, they will gain nothing. Otherwise, a pair of man and woman showing up at the same site will gain  $\alpha$  or  $\beta$  ( $\alpha > \beta$ ), i.e., when both choose *Boxing*, the man agent gains  $\alpha$  and the woman agent gains  $\beta$ ; when both choose *Ballet*, the man agent gains  $\beta$  and the woman agent gains  $\alpha$ . The population profile is denoted by  $x = (x(1), x(2), x(3), x(4))$  where  $x(1)$  is the proportion of men who chooses *Boxing*,  $x(2)$  is the proportion of men who chooses *Ballet*, and so on. The sum of  $x(1)$  and  $x(2)$  is the proportion of men, denoted by  $\omega$ , and the sum of  $x_3$  and  $x_4$  is women, which is  $1 - \omega$ . We assume that the desired equilibrium (attractor) is  $x^* = (\omega, 0, 1 - \omega, 0)$  when both man and woman choose *Boxing*. The non-equilibrium strategy of players is *Ballet*. Based on the simplex analysis, the robustness of the game should be  $((1 - \omega)\alpha)/[(1 - \omega)\alpha + \beta]$ . As one can observe in Fig. 5, the simulated robustness converges to the theoretical results as the population size is large enough.

## 6 Evaluate incentive mechanisms in reputation systems for E-marketplaces

We now apply the simulation framework to evaluate and compare the robustness of four incentive mechanisms in reputation systems for e-marketplaces to promote truthful ratings from buyers. Specifically, Jurca [13] proposes a *side-payment mechanism* where truthfully providing ratings is the buyers' equilibrium strategy. Jurca further proposes another version of this mechanism to minimize the budget required to reward buyers, called *min-budget side-payment mechanism* [13]. The *credibility mechanism* [19] requires both buyers and sellers to submit ratings about the outcome of their transactions. If there is a disagreement, both of them are punished and prevented from conducting transactions for several periods, because such a case signals that one of them is lying. The *trust-based mechanism* [36] is built on a buyer social network where reputable buyers (advisors) are more likely to be chosen as advisors by many other buyers. Sellers offer more attractive products or lower prices to satisfy those reputable buyers who will be helpful in propagating the sellers' good reputation. Buyers are thus incentivized to provide truthful ratings so as to become reputable. These reputation systems assist a buyer agent in selecting its interaction seller agents in e-marketplaces by impacting the future expected utility gain of the seller agents. It indicates that the game representing these reputation mechanisms should be in a repeated manner, which cannot be represented by payoff matrices, showing the wide applicability of the proposed simulation framework.

In the simulation framework, a set of buyers and sellers, in the proportion 10 : 1, periodically interact. In each period, each buyer chooses one seller to conduct a transaction by considering the sellers' reputation, i.e., the probability of a seller being chosen is in proportion to her reputation value. After the transaction, the buyer (and/or seller) has a chance to rate the transaction. The provided ratings are collected and processed to update seller trustworthiness (and/or buyer social network). At the end of each period, the population then evolves according to Algorithm 1, i.e., we randomly select a buyer to choose her strategy based on the expected payoff of each existing strategy and also randomly choose another buyer to mutate towards one of the extinct strategies according to Eq. (4). The strategy set of a rational buyer contains two strategies: desired equilibrium strategy and non-equilibrium strategy.<sup>7</sup> The parameters of the four incentive mechanisms are set according to their suggested values

<sup>7</sup> The strategy of rational buyers can be various, and we assume that the rational buyers would not investigate the strategies not existing in the system.

as in the corresponding papers [13, 19, 36]. In the trust-based incentive mechanism [36], the maximum number of advisors for a buyer is set to 1/4 of the whole buyer population.

We simulate four non-equilibrium strategies (i.e., attacks) adopted by irrational agents, which have been well recognized in the literature of reputation systems for e-marketplaces [13, 14]: (1) *random attack*,  $y_1$ , where a buyer randomly reports ratings for the transactions with sellers; (2) *constant attack*,  $y_2$ , where a buyer always reports the opposite ratings for transactions; (3) *whitewashing attack*,  $y_3$ , where a buyer first provides untruthful ratings, then leaves the system but re-enters using a new identity; (4) *collusive attack*,  $y_4$ , where a set of buyers form a group to conduct transactions with trustworthy sellers (top 50 %) and report negative ratings for them.<sup>8</sup> To measure the robustness of a mechanism against a specific attack, each time the simulation framework takes into account two strategies, i.e., the desired equilibrium strategy (truthfully reporting strategy)  $x = (1, 0)$  and the simulated attacking strategy  $y_k = (0, 1)(k \in \{1, 2, 3, 4\})$ . In the simulated scenarios, the population size,  $N$ , is equal to the number of rational buyers who take the desired strategy  $x$ , starting from 1. The irrational buyers taking  $y_k$ , in a number of  $M$  starting from 1, enter into the simulated system. The combined population then evolve according to the simulation framework (Algorithm 1).

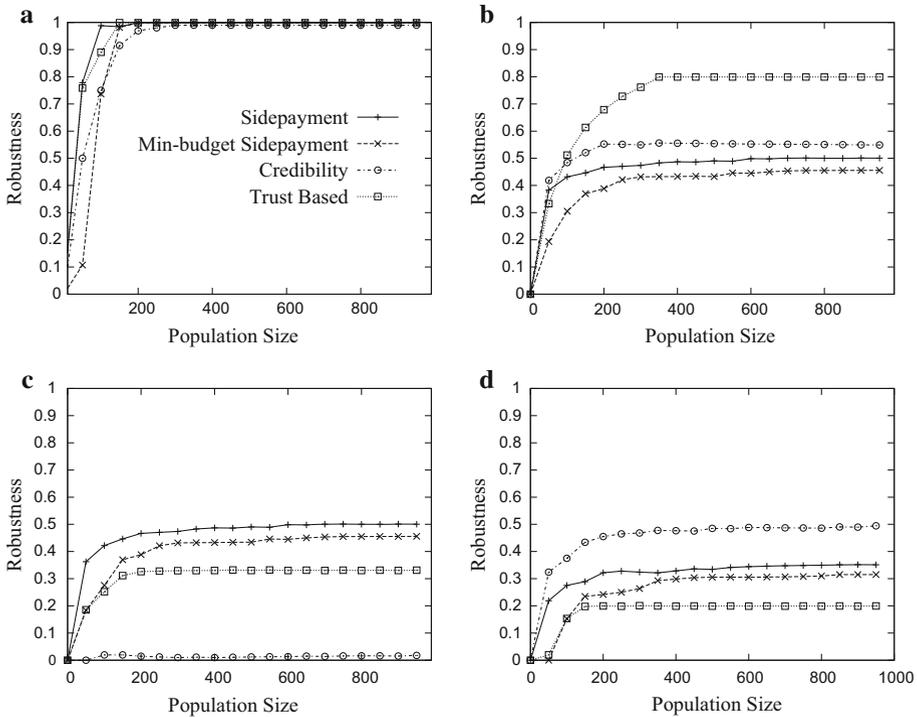
The robustness results are shown in Fig. 6. Figure 6a shows that all four mechanisms are totally robust against agents that report randomly. As the proportion of these attackers increases, the utility of buyers in taking any single non-equilibrium strategy cannot increase due to their random ratings.

Figure 6b shows the estimated robustness against the constant attack. The trust-based mechanism performs the best and min-budget side-payment mechanism performs the worst. In the trust-based mechanism, buyers choose other honest buyers (advisors) to join their social network. The constant attackers have little chance to enter rational buyers' network (i.e., gain high reputation), thus cannot significantly affect the behavior of rational buyers. In the credibility mechanism, as the proportion of the constant attackers increases, rational sellers would first deviate from equilibrium strategy, which further causes rational buyers to deviate. Thus, its robustness is only slightly higher than 0.5. In the side-payment mechanism, ratings provided by buyers are compared with the majority of buyers' ratings. When the majority becomes dishonest, the rational buyers will also deviate. Its robustness is thus around 0.5. The min-budget side-payment mechanism is less robust than the side-payment mechanism, because the former sacrifices the robustness property in order to achieve the minimal side-payment imposed on the e-marketplace owner.

Figure 6c shows the robustness against the whitewashing attack. The two side-payment mechanisms achieve higher robustness than the other two mechanisms, and the credibility mechanism bears the lowest robustness. The side-payment mechanisms have little reliance on the experience or history of buyers, resulting in similar robustness as that against the constant attack in Fig. 6b. Further, in the credibility mechanism, if a punished agent, that cannot conduct transactions, re-enters the system with a new identity, the punishment becomes invalid and fails to effectively incentivize agents to be honest, resulting in the lowest robustness.

Figure 6d shows robustness against the collusive attack. The credibility mechanism achieves the highest robustness but the trust-based mechanism performs the worst. The cred-

<sup>8</sup> The main purpose of conducting the empirical evaluation is to verify whether the proposed simulation framework can output meaningful robustness evaluation results of different incentive mechanisms against different attacks. By considering those typical pure attacking strategies, we can (1) leverage the analytical results in the literature such as [14]; (2) perform qualitative analysis on the robustness of the incentive mechanisms. Then, we can easily check whether our simulation framework implemented in the empirical evaluation gives the same results.



**Fig. 6** The robustness of the incentive mechanisms against different attacks

ibility mechanism punishes agents by not allowing them to conduct transactions for several periods, which effectively decreases the impact of the collusive attackers. In the trust-based mechanism, collusive attackers can provide untruthful ratings, add each other in the social network, and become reputable. These “reputable” attackers will then gain a high utility, causing rational buyers to deviate from the equilibrium strategy. From Fig. 6, we can observe that the robustness decreases as attacks become more refined, e.g., whitewashing and collusive attacks are more sophisticated than random and constant attacks. This trend was also predicted by other researchers [13]. In addition, none of the incentive mechanisms can always perform perfectly against different attacks.

### 7 Conclusion and discussion

In this article, we propose a quantitative approach to evaluate the robustness of incentive mechanisms against bounded rational strategies. The main contributions are summarized as follows: (1) We provide formal definitions of a robustness measure for incentive mechanisms in the presence of bounded rational players. (2) We propose a general simulation framework based on evolutionary dynamics, and analyze the influence of its parameter settings on the robustness measure introduced. (3) We then validate the outcomes of our simulation framework by comparing with theoretical predictions of robustness in several representative game models. (4) We use our simulation framework to evaluate and compare the robustness of several incentive mechanisms proposed for reputation systems in e-marketplaces against dif-

ferent non-equilibrium strategies; results show that the robustness of mechanisms decreases against sophisticated non-equilibrium strategies and highlights the need for more robust incentive mechanisms for reputation systems in e-marketplaces.

It is worth noting that the proposed robustness evaluation framework only considers a single type of non-equilibrium (non-cooperative) strategy each time, which is due to the following two reasons. Firstly, in some domains it is a typical way to separately evaluate each type of attacks, such as e-marketplaces [9, 11], even though various types of attacks could coexist in a realistic environment. Secondly, the proposed robustness is based on the concept of evolutionary stable strategy in evolutionary game theory, where the alternative strategies in consideration are generally referred to as pure strategies in non-cooperative interactions by default [23, 24, 32]. In future work, one nature extension is to consider various or mixed types of attacks (including coordinate non-equilibrium strategy) through searching for the support from other theoretical concepts other than ESS.

Moreover, the proposed simulation framework can be enhanced in several aspects. First, it can be used to search for the worst non-equilibrium strategy. In the proposed simulation framework, we are to find the robustness of an incentive mechanism against a specific non-equilibrium strategy. It is a natural extension to search for the worst attacking strategy where the incentive mechanism under evaluation exhibits the lowest robustness value. It then becomes an optimization problem to minimize the robustness in a space composed by all possible non-equilibrium strategies. Secondly, the robustness evaluation framework can be extended to allow rational players to form coalitions and conduct coordinate attacks. A coalition is formed by  $k$  rational players who coordinate their actions to increase the overall utility of the coalition, assuming that the players will use side-payments within the formed coalition to share corresponding gains. The rational coalition is different from the collusive attack used in Sect. 6 where the bounded rational players form coalitions to attack the system for gaining more utility. By checking the evolution of rational non-collusive players, our framework can be extended to calculate the largest collusive player size such that non-collusive players still sustain equilibrium strategies. Finally, the robustness of various incentive mechanisms in e-marketplaces will be evaluated through the proposed simulation framework. Using the data from these experiments, we will analyze the key factors that impact the robustness of incentive mechanisms. Our aim in this process is to design and develop more robust incentive mechanisms by carefully considering those influential factors.

## References

1. Aghassi, M., & Bertsimas, D. (2006). Robust game theory. *Mathematical Programming: Series A and B*, 107(1–2), 231–273.
2. Aiyer, A. S., Alvisi, L., & Clement, A. (2005). Bar fault tolerance for cooperative services. In *Proceedings of ACM symposium on operating systems principles (SOSP)*, pp. 45–58.
3. Bergemann, D., & Morris, S. (2012). *Robust mechanism design: The role of private information and higher order beliefs*. Singapore and London: World Scientific Publishing and Imperial College Press.
4. Díaz, J., Goldberg, L. A., Mertziou, G. B., Richerby, D., Serna, M., & Spirakis, P. G. (2012). Approximating fixation probabilities in the generalized moran process. In *Proceedings of the ACM-SIAM symposium on discrete algorithms (SODA)*, pp. 954–960.
5. Eliaz, K. (2002). Fault tolerant implementation. *Review of Economic Studies*, 69(3), 589–610.
6. Halpern, J. Y. (2008). Beyond Nash equilibrium: Solution concepts for the 21st century. In *Proceedings of the 27th ACM symposium on principles of distributed computing (PODC)*, pp. 1–10.
7. Halpern, J. Y., Pass, R., & Seeman, L. (2012). I'm doing as well as i can: Modelling people as rational finite automata. In *Proceedings of the AAAI conference on artificial intelligence*, pp. 1917–1923.
8. Halbing, D., & Yu, W. (2009). The outbreak of cooperation among success-driven individuals under noisy conditions. In *Proceedings of the National Academy of Sciences*, pp. 3680–3685.

9. Irissappane, A. A., Jiang, S., & Zhang, J. (2012). Towards a comprehensive testbed to evaluate the robustness of reputation systems against unfair rating attacks. In *Proceedings of the twentieth conference on user modeling, adaptation, and personalization (UMAP) workshop on trust, reputation, and user modelling*.
10. Jackson, M. O. (2003). *Mechanism design*. Oxford, UK: Encyclopedia of Life Support Systems (EOLSS) Publisher.
11. Jiang, S. (2013). Towards the design of robust trust and reputation systems. In *Proceedings of the twenty-third international conference on artificial intelligence (IJCAI)*, pp. 3225–3226.
12. Jøsang, A., Ismail, R., & Boyd, C. (2007). A survey of trust and reputation systems for online service provision. *Decision Support System*, 43(2), 795–825.
13. Jurca, R. (2007). Truthful reputation mechanisms for online systems. Ph.D. thesis, EPFL.
14. Kerr, R., & Cohen, R. (2009). Smart cheaters do prosper: Defeating trust and reputation systems. In *AAMAS*, pp. 993–1000.
15. Leyton-Brown, K. (2013). *Chapter 7: Mechanism design and auctions in multiagent systems* (2nd ed.). Cambridge, MA: MIT Press.
16. Moran, P. A. P. (1962). *The statistical process of evolutionary theory*. Oxford, UK: Clarendon Press.
17. Nowak, M. A. (2006). *Evolutionary dynamics: Exploring the equations of life*. Cambridge: Harvard University Press.
18. Nowak, M. A., Sasaki, A., Taylor, C., & Fudenberg, D. (2004). Emergence of cooperation and evolutionary stability in finite populations. *Nature*, 428, 646–650.
19. Papaioannou, T. G., & Stamoulis, G. D. (2010). A mechanism that provides incentives for truthful feedback in p2p systems. *Electronic Commerce Research*, 10(3–4), 331–362.
20. Ponsen, M., Tuyls, K., Kaisers, M., & Ramon, J. (2009). An evolutionary game-theoretic analysis of poker strategies. *Entertainment Computing*, 1(1), 39–45.
21. Raghunandan, M. A., & Subramanian, C. A. (2012). Sustaining cooperation on networks: An analytical study based on evolutionary game theory. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pp. 913–920.
22. Sandholm, W. H. (2012). Evolutionary game theory. *Computational Complexity*, 1, 1000–1029.
23. Smith, J. M. (1982). *Evolution and the theory of games*. Cambridge: Cambridge University Press.
24. Smith, J. M., & Price, G. R. (1973). The logic of animal conflict. *Nature*, 246(5427), 15–18.
25. Taylor, P. D. (1978). Evolutionarily stable strategies and game dynamics. *Mathematical Biosciences*, 6, 145–156.
26. Traulsen, A., Hauert, C., Silva, H. D., Nowak, M. A., & Sigmund, K. (2009). Exploration dynamics in evolutionary games. In *Proceedings of National Academy of Sciences of the United States of America*, pp. 709–712.
27. Tuyls, K., & Parsons, S. (2007). What evolutionary game theory tells us about multiagent learning. *Artificial Intelligence*, 171(7), 406–416.
28. Tuyls, K., Hoen, P. J., & Vanschoenwinkel, B. (2006). An evolutionary dynamical analysis of multi-agent learning in iterated games. *Journal of Autonomous Agents and Multi-Agent Systems*, 12(1), 115–153.
29. Ventura, D. (2012). Rational irrationality. In *Proceedings of the AAAI spring symposium on game theory for security sustainability and health*, pp. 83–90.
30. von Neumann, J., & Morgenstern, O. (1944). *Theory of games and economic behaviour*. Princeton: Princeton University Press.
31. Walsh, W. E., Das, R., Tesauro, G., & Kephart, J. O. (2002). Analyzing complex strategic interactions in multi-agent systems. In *In AAAI-03 workshop on game theoretic and decision theoretic agents*, pp. 109–118.
32. Wellman, M. P. (1997). *Evolutionary game theory*. Cambridge: MIT Press.
33. Wellman, M. P. (2006). Methods for empirical game-theoretic analysis. In *Proceedings of the twenty-first national conference on artificial intelligence (AAAI)*, pp. 1552–1553.
34. Witkowski, J., Seuken, S., & Parkes, D. C. (2011). Incentive-compatible escrow mechanisms. In *Proceedings of the 25th conference on artificial intelligence (AAAI)*, pp. 751–757.
35. Wright, J. R., & Leyton-Brown, K. (2013). Behavioral game-theoretic models: A Bayesian framework for parameter analysis. In *Proceedings of the 11th international conference on autonomous agents and multiagent systems (AAMAS)*, pp. 921–928.
36. Zhang, J., Cohen, R., & Larson, K. (2012). Combining trust modeling and mechanism design for promoting honesty in e-marketplaces. *Computational Intelligence*, 28(4), 549–578.