

# To Trust or Not to Trust: Evolutionary Dynamics of an Asymmetric N-player Trust Game

Ik Soo Lim and Naoki Masuda

**Abstract**—Trusting others and reciprocating the received trust with trustworthy actions are fundamentals of economic and social interactions. The trust game (TG) is widely used for studying trust and trustworthiness and entails a sequential interaction between two players, an investor and a trustee. It requires at least two strategies or options for an investor (e.g. to trust versus not to trust a trustee). According to the evolutionary game theory, the antisocial strategies (e.g. not to trust) evolve such that the investor and trustee end up with lower payoffs than those that they would get with the prosocial strategies (e.g. to trust). A generalisation of the TG to a multiplayer (i.e. more than two players) TG was recently proposed. However, its outcomes hinge upon two assumptions that various real situations may substantially deviate from: (i) investors are forced to trust trustees and (ii) investors can turn into trustees by imitation and vice versa. We propose an asymmetric multiplayer TG that allows investors not to trust and prohibits the imitation between players of different roles; instead, investors learn from other investors and the same for trustees. We show that the evolutionary game dynamics of the proposed TG qualitatively depends on the nonlinearity of the payoff function and the amount of incentives collected from and distributed to players through an institution. We also show that incentives given to trustees can be useful and sufficient to cost-effectively promote trust and trustworthiness among self-interested players.

**Index Terms**—Evolutionary game theory, evolutionary dynamics, replicator dynamics, trust game, incentives

## I. INTRODUCTION

The evolution of pro-social behaviours among self-interested individuals has been a focus of research across disciplines. For instance, the evolution of cooperation in social dilemma situations such as the Prisoner’s Dilemma (PD) and its  $N$ -player generalisation, the Public Goods Game (PGG), has attracted lots of attention [1][2][3][4][5][6]. Evolutionary game theory provides a theoretical framework with which to study the evolution of strategies or behaviours among self-interested individuals in these social dilemmas or other situations, in which successful strategies or genes are spread by fitness-dependent reproduction and imitation [7][8]. It has also been widely used for applications such as modelling the propagation of competing technologies and policies for green supply chain management [9][10].

Non-simultaneous or sequential interactions between two players are common in many situations such as buyer-seller interactions, whereas the PD and PGG are concerned with simultaneous interactions. Non-simultaneous interactions yield a problem of trust in the sense that the decision by one of two

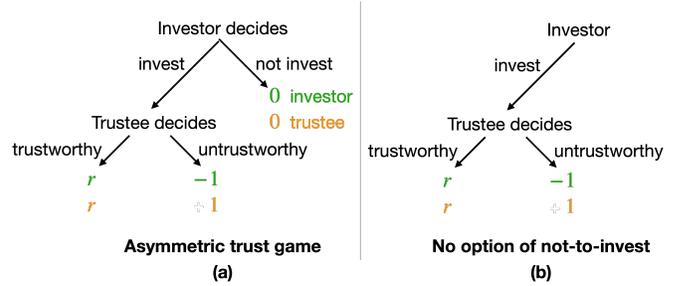


Fig. 1. Two-player binary TGs. (a) Game tree of the asymmetric two-player binary TG, referred to as a general two-player TGIG in Ref. [27], in which the role of each player is fixed. The payoffs of an investor are shown in green. Those of a trustee are shown in orange. Adapted from Ref. [19]. We generalise this game to an  $N$ -player game in this article. (b) Game tree of the two-player binary TG that is used for the generalisation to the NTG in Ref. [27]. This game does not allow an investor not to invest. In both (a) and (b), we require  $0 < r < 1$ , where  $r$  represents the relative productivity of the prosocial strategies.

players (e.g. a buyer) can make oneself vulnerable to potential exploitation by the other (e.g. a seller) [11]. In such situations, higher levels of trusting in others and reciprocating the received trust with trustworthy actions have been associated with more efficient judicial systems, higher quality in government bureaucracies, lower corruption, greater financial development, and better economic outcomes among other benefits for the society [12]. The concept of trust has also attracted interest in engineering research communities, ranging from networking to human-machine interaction and artificial intelligence [13][14][15], where many problems are cast as buyer-seller interactions [16]. The trust game (TG) is a current gold standard of formalisation for non-simultaneous interaction in social dilemma situations and has widely been used to study trust and trustworthiness [11][12][17][18][19][20][21][22][23][24]. The TG is composed of a one-shot sequential interaction between two players in different roles, one as an investor (representing, for example, a truster, buyer, or citizen) and the other as a trustee (representing, for example, a seller or governor). One of the simplest variants of TGs is the binary TG, which involves two strategies per role [19][25][26]. An investor either invests (i.e. trusts) or does not invest in a trustee. Then, the trustee decides to be either trustworthy or untrustworthy to the investor (Fig. 1a).

The evolutionary game theory predicts that self-interested strategies (e.g. for an investor not to invest) evolve in the two-player binary TG. The classical game theory also yields a similar conclusion via backward induction; given investment from an investor, a rational trustee is better off by being untrustworthy and, anticipating it, a rational investor does not

Ik Soo Lim is with School of Computing and Mathematical Sciences, University of Greenwich, London SE10 9LS, U.K. (e-mail: i.lim@gre.ac.uk) Naoki Masuda is with Department of Mathematics and Computational and Data-Enabled Science and Engineering Program, State University of New York at Buffalo, USA

Manuscript received XXXX; revised YYYY.

invest in a trustee in the first place. Thus, the two players end up with lower payoffs than those that they would get with the pro-social strategies (i.e. for the investor to invest and for the trustee to be trustworthy). Therefore, an additional mechanism is required for promoting the evolution of the pro-social strategies in the TG [28][29][30].

An  $N$ -player binary TG (NTG) was recently proposed as a multiplayer (i.e.  $N \geq 2$ ) generalisation of the binary TG [27]. However, it suffers from two major difficulties that hamper us from clarifying mechanisms of trust and trustworthiness in multiplayer situations in reasonably realistic manners. First, in this NTG, the investor does not have an option not to invest (Fig. 1b); the investor is assumed to invest. Therefore, one cannot investigate the evolution and stability of trusting as opposed to non-trusting behaviour. Note that their NTG with  $N = 2$  players is not the two-player TG, which this model attempted to generalise. Second, investors are allowed to turn into trustees and vice versa by payoff-driven imitation. An evolutionary outcome of this second assumption is the cease of game playing because all players eventually become trustees [27]. Without an investor, one cannot carry on the game. The justification of this result and the underlying assumption of the role-unaware imitation is unclear. The NTG with citizens and governors was used as an example in Ref. [27], where citizens were allowed to imitate and become governors. The evolutionary outcome is that all players become governors. Once there is no citizen, there is no NTG to be played. A population composed of all governors but no citizen is not only unrealistic but also incompatible with the behavioural experiment setups of the TGs, which ensures that both a citizen (or an investor) and a governor (or a trustee) are always available to play the TG [12]. The follow-up studies of the original NTG [27] also inherit the aforementioned two assumptions, i.e., that the investor does not have a choice not to invest and that players can turn into a preferred role by imitation [31][32][33].

In reality, investor-trustee interactions often involve multiplayer interactions rather than dyadic ones; for instance, multiple investors may be involved in a large project. Hence, setting up reasonable NTGs and understanding their population dynamics remains a worthwhile goal. Our contributions in this paper are threefold:

- We propose an asymmetric NTG with two strategies per role, which generalises the two-player TG but does not suffer from the two problems inherent in the previously proposed NTG.
- We introduce non-linear payoff functions that can yield evolutionary dynamics qualitatively different from that of a linear one.
- We propose an incentive scheme to cost-effectively steer the self-interested players to take prosocial strategies such that the population average of the payoff (or social welfare) is maximised.

The source code used for this paper is provided on Github: [https://github.com/iksoolim/asymmetric\\_N-player\\_trust\\_game](https://github.com/iksoolim/asymmetric_N-player_trust_game).

## II. MODEL

### A. Population and Group Formation

We consider an asymmetric NTG in which the role of each individual is fixed as either investor or trustee throughout the whole evolutionary dynamics. Furthermore, we assume that social learning, i.e., payoff-led imitation of strategies, only occurs among individuals of the same role as in the two-player TG [19]. There are two strategies available for each role. An investor either invests or does not invest in trustees. A trustee selects to be either trustworthy or untrustworthy to investors. We consider two infinitely large populations, one for investors and the other for trustees. From time to time, a group of  $N_I$  investors and  $N_T$  trustees, selected uniformly at random from the respective population, is formed and these  $N \equiv N_I + N_T$  individuals participate in a one-shot NTG. We assume that  $N_I$  and  $N_T$  are fixed.

### B. Payoffs

We assume that the total value of the investment aggregated over the investing investors is equal to

$$\frac{1 - w^{k_i}}{1 - w} = \begin{cases} 0 & \text{if } k_i = 0, \\ 1 & \text{if } k_i = 1, \\ 1 + w + w^2 + \dots + w^{k_i-1} & \text{if } k_i \geq 2, \end{cases} \quad (1)$$

where  $k_i \in \{0, 1, \dots, N_I\}$  denotes the number of investing investors in the group, and  $w > 0$  determines how the value of the investments accumulates when an additional investor contributes to the collective good. A similar non-linear payoff function was previously used for the PGG [4]. If  $0 < w < 1$ , then the value of the contribution by each additional investing investor is diminishing, i.e. discounted or sub-additive. If  $w = 1$ , then the value of the contribution is 1 for any investor regardless of the number of investing investors,  $k_i$ . This linear payoff function is the same as that for the original NTG [27]. Note that the total value of the investment is equal to  $k_i$  when  $w = 1$ , which follows from L'Hopital's rule applied to the left-hand side of Eq.(1). If  $w > 1$ , the value of the contribution per investor increases as  $k_i$  increases, i.e. representing synergistic or super-additive benefits.

The total investment is equally divided and distributed to the  $N_T$  trustees. Therefore, the payoff that an untrustworthy trustee in the group receives from the game, denoted by  $\Pi_u^o(k_i)$ , is given by

$$\Pi_u^o(k_i) = \frac{1}{N_T} \frac{1 - w^{k_i}}{1 - w}. \quad (2)$$

The payoff of a trustworthy trustee in the group, denoted by  $\Pi_t^o(k_i)$ , is given by

$$\Pi_t^o(k_i) = r \Pi_u^o(k_i) = r \frac{1}{N_T} \frac{1 - w^{k_i}}{1 - w}, \quad (3)$$

where  $r$  represents relative productivity of the prosocial strategies and satisfies  $0 < r < 1$ . In the two-player TG, when an investing investor and a trustworthy trustee interact with each other, each of them gets the same payoff (Fig. 1a). In the

$N$ -player generalisation, analogously, we assume that when a group of investing investors and a group of trustworthy trustees interact with each other, each group gets the same (group) payoff. The aggregated return from the  $k_t \in \{0, 1, \dots, N_T\}$  trustworthy trustees is equally distributed to the  $k_i$  investing investors in the group. Therefore, the payoff that an investing investor receives from the game, denoted by  $\Pi_i^o(k_i, k_t)$ , is given by

$$\begin{aligned} \Pi_i^o(k_i, k_t) &= \underbrace{\frac{1}{k_i} k_t \Pi_t^o(k_i)}_{\text{net gain}} + \underbrace{(N_T - k_t) \left(-\frac{1}{N_T}\right)}_{\text{net loss}} \\ &= \frac{k_t}{N_T} \frac{r(1-w^{k_i})}{k_i(1-w)} + \left(1 - \frac{k_t}{N_T}\right) \cdot (-1). \end{aligned} \quad (4)$$

The payoff  $\Pi_i^o(k_i, k_t)$  is equal to the expected payoff of an investing investor playing a two-player game with each of the  $N_T$  trustees; the net gain from a trustworthy trustee is  $\frac{r(1-w^{k_i})}{k_i(1-w)}$  and the net loss from an untrustworthy trustee is  $-1$ . Lastly, the payoff of a non-investing investor is  $\Pi_n^o = 0$ . Note that a special case of  $N_I = N_T = 1$  recovers the two-player TG (Fig. 1a).

By including incentives and associated costs for the players, we define the final payoffs  $\Pi_i$ ,  $\Pi_n$ ,  $\Pi_t$ , and  $\Pi_u$  for an investing investor, non-investing investor, trustworthy trustee and untrustworthy trustee, respectively, by

$$\Pi_i(k_i, k_t) = \Pi_i^o(k_i, k_t) + v_I - av_I, \quad (5)$$

$$\Pi_n = \Pi_n^o - av_I, \quad (6)$$

$$\Pi_t(k_i) = \Pi_t^o(k_i) + v_T - av_T, \quad (7)$$

$$\Pi_u(k_i) = \Pi_u^o(k_i) - av_T, \quad (8)$$

where an investor pays a fee  $av_I$  to the institution providing the incentives and an investing investor receives a reward  $v_I$  from the institution, where  $v_I \geq 0$ . We assume the fee rate  $a > 1$  such that the total incentive is less than the total fee, taking into consideration the operating cost for the institution. Similarly, a trustee pays a fee  $av_T$  to the institution and a trustworthy trustee receives a reward  $v_T \geq 0$ . A similar incentive scheme has been assumed for the PGG [6]. For a given investor in a group of  $N$  players, the probability that  $m_t$  among  $N_T$  trustees are trustworthy (and thus  $N_T - m_t$  trustees are untrustworthy) is  $\binom{N_T}{m_t} y_t^{m_t} (1-y_t)^{N_T-m_t}$ , where  $y_t$  denotes the fraction of trustworthy trustees in the trustee population;  $1-y_t$  is the fraction of untrustworthy trustees. For a given investor, the probability that  $m_i$  among the other  $N_I-1$  investors are investing is  $\binom{N_I-1}{m_i} y_i^{m_i} (1-y_i)^{N_I-1-m_i}$ , where  $y_i$  denotes the fraction of investing investors in the investor population. Therefore, the expected payoff for an investing investor is

$$\begin{aligned} P_i &= \sum_{m_i=0}^{N_I-1} \binom{N_I-1}{m_i} y_i^{m_i} (1-y_i)^{N_I-1-m_i} \\ &\quad \times \sum_{m_t=0}^{N_T} \binom{N_T}{m_t} y_t^{m_t} (1-y_t)^{N_T-m_t} \Pi_i(m_i+1, m_t) \\ &= \frac{r}{N_I(1-w)} \frac{y_t}{y_i} \left\{ 1 - [1 + (w-1)y_i]^{N_I} \right\} + y_t - 1 \\ &\quad + v_I - av_I. \end{aligned} \quad (9)$$

Similarly, the expected payoffs  $P_n$ ,  $P_t$  and  $P_u$  for a non-investing investor, trustworthy trustee and untrustworthy trustee, respectively, are given by

$$P_n = -av_I, \quad (10)$$

$$\begin{aligned} P_t &= \sum_{m_i=0}^{N_I} \binom{N_I}{m_i} y_i^{m_i} (1-y_i)^{N_I-m_i} \\ &\quad \times \sum_{m_t=0}^{N_T-1} \binom{N_T-1}{m_t} y_t^{m_t} (1-y_t)^{N_T-1-m_t} \Pi_t(m_i), \\ &= \frac{r}{N_T(1-w)} \left\{ 1 - [1 + (w-1)y_i]^{N_I} \right\} + v_T - av_T, \end{aligned} \quad (11)$$

$$P_u = \frac{1}{N_T(1-w)} \left\{ 1 - [1 + (w-1)y_i]^{N_I} \right\} - av_T. \quad (12)$$

See Appendix A for the derivation Eqs. (9), (11) and (12).

### C. Evolutionary Game Dynamics

For the evolutionary game dynamics, we use asymmetric replicator equations given by

$$\begin{aligned} \dot{y}_i &= y_i(P_i - P_I) = y_i(1-y_i)(P_i - P_n) \\ &= (1-y_i)y_i \left( \frac{ry_t \left\{ 1 - [1 + (w-1)y_i]^{N_I} \right\}}{N_I(1-w)y_i} + y_t - 1 + v_I \right), \end{aligned} \quad (13)$$

$$\begin{aligned} \dot{y}_t &= y_t(P_t - P_T) = y_t(1-y_t)(P_t - P_u) \\ &= (1-y_t)y_t \left( \frac{(r-1) \left\{ 1 - [1 + (w-1)y_i]^{N_I} \right\}}{N_T(1-w)} + v_T \right), \end{aligned} \quad (14)$$

where the dot denotes a time derivative,  $P_I = y_i P_i + (1-y_i) P_n$  is the average payoff of the investor in the entire population, and  $P_T = y_t P_t + (1-y_t) P_u$  is the average payoff of the trustee.

To analyse the dynamics given by Eqs. (13) and (14), we find all equilibria by setting  $\dot{y}_i = \dot{y}_t = 0$ . The stability of an equilibrium is determined by the eigenvalues of the Jacobian matrix, which is given by

$$J = \begin{pmatrix} \frac{\partial \dot{y}_i}{\partial y_i} & \frac{\partial \dot{y}_i}{\partial y_t} \\ \frac{\partial \dot{y}_t}{\partial y_i} & \frac{\partial \dot{y}_t}{\partial y_t} \end{pmatrix} = \begin{pmatrix} J_{11} & J_{12} \\ J_{21} & J_{22} \end{pmatrix} \quad (15)$$

at the equilibrium, where

$$\begin{aligned} J_{11} &= r(1-y_i)y_t \left[ (w-1)y_i + 1 \right]^{N_I-1} \\ &\quad - \frac{ry_t \left\{ [(w-1)y_i + 1]^{N_I} - 1 \right\}}{N_I(w-1)} - (2y_i-1)(v_I + y_t - 1), \end{aligned} \quad (16)$$

$$J_{12} = \frac{(1-y_i) \left( r \left\{ [(w-1)y_i + 1]^{N_I} - 1 \right\} + N_I(w-1)y_i \right)}{N_I(w-1)}, \quad (17)$$

$$J_{21} = \frac{N_I(r-1)(1-y_t)y_t \left[ (w-1)y_i + 1 \right]^{N_I-1}}{N_T}, \quad (18)$$

$$J_{22} = \frac{(2y_t-1)(r-1) \left\{ 1 - [(w-1)y_i + 1]^{N_I} \right\}}{N_T(w-1)}$$

$$-(2y_t - 1)v_T. \quad (19)$$

If any of the two eigenvalues is positive, the equilibrium is unstable. Otherwise, the equilibrium is stable; trajectories starting close enough to the equilibrium remain close enough. Especially, the equilibrium is asymptotically stable if and only if all the eigenvalues are negative; in this case, trajectories starting close enough to the equilibrium converge to it [34]. Note that Eqs. (13), (14) (16), (17) and (19) are also valid for  $w = 1$  with the use of L'Hopital's rule.

### III. RESULTS

In this section, we characterize the equilibria, their stability, and trajectories of the dynamical system given by Eqs. (13) and (14), of which the state space is  $\{(y_i, y_t) \in [0, 1]^2\}$ . Note that Eqs. (13) and (14) imply that  $(0, 0)$ ,  $(0, 1)$ ,  $(1, 0)$ , and  $(1, 1)$  are always equilibria. For proof of the stability of these and the other equilibria, see Appendix B.

#### A. $v_T = 0$

For  $v_T = 0$ , the edge  $y_i = 0$  of the state space is a line of equilibria. For  $v_T = 0 \wedge 0 \leq v_I < 1$ , the part of the edge satisfying  $0 \leq y_t < \frac{1-v_I}{r+1}$ , including the origin,  $(y_i, y_t) = (0, 0)$ , is stable but not asymptotically stable (Fig. 2a). The points on the line satisfying  $\frac{1-v_I}{r+1} < y_t \leq 1$ , including  $(0, 1)$ , as well as  $(1, 0)$  and  $(1, 1)$ , are unstable equilibria. As Fig. 2a indicates, any trajectory is eventually attracted to one of the stable equilibria. This evolutionary outcome is qualitatively the same as that of the two-player TG and it is so irrespectively of the non-linearity  $w$  in the payoff function (e.g. for any of  $w \in \{0.6, 1, 1.4\}$ ). With the special case of  $v_T = 0 \wedge v_I = 0 \wedge w = 1$ , we obtain a baseline model, which is an  $N$ -player generalisation of the two-player TG without any other mechanism.

For  $v_T = 0 \wedge v_I > 1$ , the equilibrium  $(1, 0)$  is not only asymptotically stable but also globally convergent (i.e. reached from any initial state). The equilibria  $(0, 0)$ ,  $(0, 1)$ ,  $(1, 1)$  and  $y_i = 0$  are unstable.

#### B. $0 < v_T < v_T^*$

For  $0 < v_T < v_T^* \equiv \frac{(1-r)(w^{N_I}-1)}{N_T(w-1)} \wedge 0 \leq v_I < 1$ , an interior equilibrium

$$\mathbf{Q} = \left( \frac{d^{1/N_I} - 1}{w - 1}, \frac{N_I(1 - v_I)(d^{1/N_I} - 1)}{N_I(d^{1/N_I} - 1) + (d - 1)r} \right), \quad (20)$$

emerges, where  $d = 1 + \frac{N_T v_T (w - 1)}{1 - r}$ . The interior equilibrium is at the intersection of the two nullclines,  $P_i - P_n = 0$  and  $P_t - P_u = 0$  with  $0 < y_i < 1 \wedge 0 < y_t < 1$ ; see Appendix B5 for the proof of the existence of the interior equilibrium. Note that L'Hopital's rule implies that  $v_T^* = \frac{N_I(1-r)}{N_T}$  and  $\mathcal{Q} = \left( \frac{N_T v_T}{N_I(1-r)}, \frac{1-v_I}{1+r} \right)$  for  $w = 1$ .

The interior equilibrium is asymptotically stable for  $w < 1$ , neutrally stable for  $w = 1$ , and unstable for  $w > 1$  (Fig. 2b). The other equilibria are the four corners of the state space, all of which are unstable. For  $w = 1$ , at which all the trajectories

surrounding  $\mathbf{Q}$  form closed cycles, the time average of  $(y_i, y_t)$  over each of the cycles is equal to  $(y_i, y_t)$  at  $\mathbf{Q}$  given by Eq. (20); see Appendix C for the proof. For  $w > 1$ , all the trajectories converge to the heteroclinic cycle consisting of the four unstable equilibria, which are saddle points, and the four edges that connect them;  $(0, 0) \rightarrow (0, 1) \rightarrow (1, 1) \rightarrow (1, 0) \rightarrow (0, 0)$ . In this case, the time average of  $y_i$  and  $y_t$  over the heteroclinic cycle does not converge; see Appendix D for the proof.

For  $0 < v_T < v_T^* \wedge v_I > 1$ , there does not exist any interior equilibrium. In this case, only the four corners are equilibria. The equilibrium  $(1, 0)$  is not only asymptotically stable but also globally convergent. The equilibria  $(0, 0)$ ,  $(0, 1)$  and  $(1, 1)$  are unstable.

#### C. $v_T = v_T^*$

At  $v_T = v_T^*$ , a line of equilibria  $y_i = 1$  emerges. For  $v_T = v_T^* \wedge 0 \leq v_I < 1$ , the part of the line satisfying  $\frac{N_I(1-v_I)(w-1)}{r(w^{N_I}-1)+N_I(w-1)} < y_t \leq 1$ , including  $(y_i, y_t) = (1, 1)$ , is stable but not asymptotically stable (Fig. 2c). The part of the line satisfying  $0 \leq \frac{N_I(1-v_I)(w-1)}{r(w^{N_I}-1)+N_I(w-1)} < y_t$ , including  $(1, 0)$ , and the equilibria  $(0, 0)$  and  $(0, 1)$  are unstable.

For  $v_T = v_T^* \wedge v_I > 1$ , the whole line of equilibria including  $(1, 0)$  and  $(1, 1)$  is stable but not asymptotically stable. The equilibria  $(0, 0)$  and  $(0, 1)$  are unstable. These results are qualitatively the same across the different  $w$  values.

#### D. $v_T > v_T^*$

For  $v_T > v_T^*$ , only the four corners are equilibria. The equilibrium  $(1, 1)$  is not only asymptotically stable but also globally convergent (Fig. 2d). Note that  $(1, 1)$  represents the fully cooperative populations entirely consisting of investing investors and trustworthy trustees. All the other equilibria, namely,  $(0, 0)$ ,  $(0, 1)$  and  $(1, 0)$ , are unstable. These results hold true independently of the  $v_I \geq 0$  and  $w$  values, except for the dependence of  $v_T^*$  on  $w$ .

In Fig. 3, we show a schematic diagram summarising the analysis so far. It presents the evolutionary dynamics that varies in a qualitatively different manner depending on the incentive values  $v_I$  and  $v_T$ .

#### E. Population Average of Payoff and Optimal Incentive

One of our goals for proposing and analysing the present NTG is to steer the self-interested players to behave pro-socially, increase the efficiency of the equilibrium in terms of the payoff the players gain and do so in a cost-efficient manner. Therefore, in this section, we analyze the population average of the payoff given by

$$\begin{aligned} P(y_i, y_t) &= \frac{N_I}{N_I + N_T} P_I(y_i, y_t) + \frac{N_T}{N_I + N_T} P_T(y_i, y_t) \\ &= -\frac{a(w-1)(N_I v_I + N_T v_T) + 1}{(w-1)(N_I + N_T)} + \frac{N_I(v_I - 1)}{N_I + N_T} y_i \\ &\quad + \frac{N_T v_T (w-1) - 2r + 1}{(w-1)(N_I + N_T)} y_t + \frac{N_I}{N_I + N_T} y_i y_t \\ &\quad + \frac{[(2r-1)y_t + 1][(w-1)y_i + 1]^{N_I}}{(w-1)(N_I + N_T)} \end{aligned} \quad (21)$$

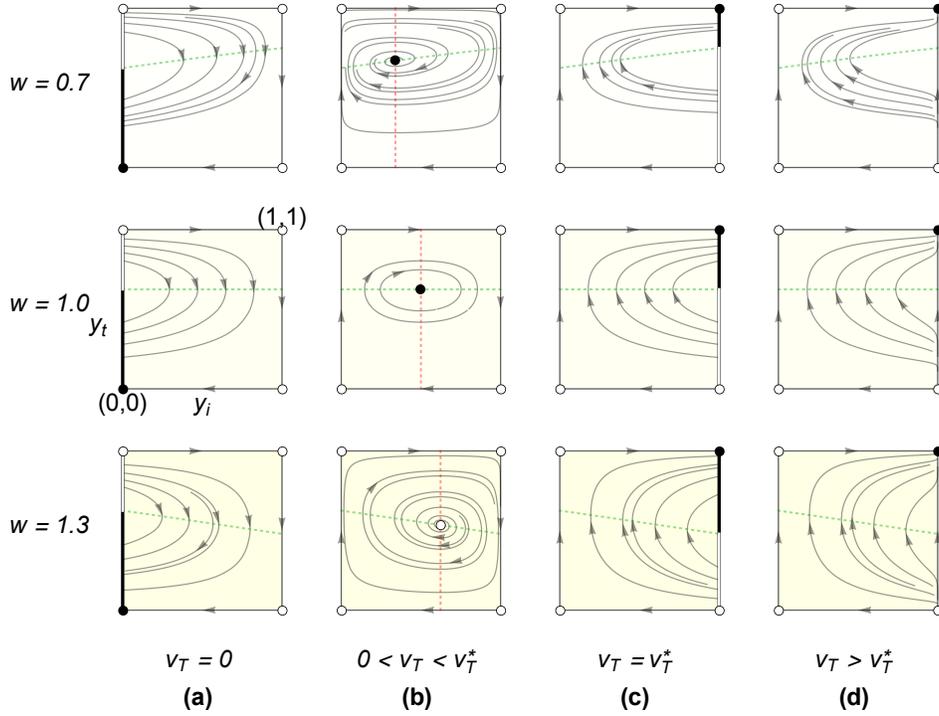


Fig. 2. Evolutionary game dynamics of the asymmetric NTG with fixed roles for the players. We set  $N_I = 5, N_T = 5, r = 0.6, v_I = 0$  and  $v_T/v_T^* \in \{0, 0.5, 1, 1.1\}$ . (1st row)  $w = 0.7$ , (2nd)  $w = 1$ , and (3rd)  $w = 1.3$ . A filled circle represents a stable equilibrium. An open circle represents an unstable equilibrium. On edges  $y_i = 0$  and  $y_t = 1$ , the thick solid lines indicate stable equilibria and the hollow lines indicate unstable equilibria. The dashed lines indicate the nullclines  $P_i - P_n = 0$  (in green) and  $P_t - P_u = 0$  (in red). (a) When  $v_T = 0$  (i.e. no incentive to trustworthy trustees), all trajectories converge to a lower part of the edge  $y_i = 0$ , and investment (i.e. trust) does not evolve. (b) When  $0 < v_T < v_T^*$ , an interior equilibrium point emerges and moves, with increasing  $v_T$ , from  $y_i = 0$  towards  $y_i = 1$ . (c) When  $v_T = v_T^*$ , the interior equilibrium disappears and all trajectories converge to an upper part of the edge  $y_i = 1$ . (d) When  $v_T > v_T^*$ , all trajectories converge to  $(1, 1)$ , i.e., the state of full trust and full trustworthiness. The nonlinearity in the payoff function yields a stable interior equilibrium with trajectories spiralling into it or an unstable interior equilibrium with trajectories spiralling out of it. These dynamics are qualitatively different from those in the case of the linear payoff function (i.e. a neutrally stable interior equilibrium with periodic trajectories around it).

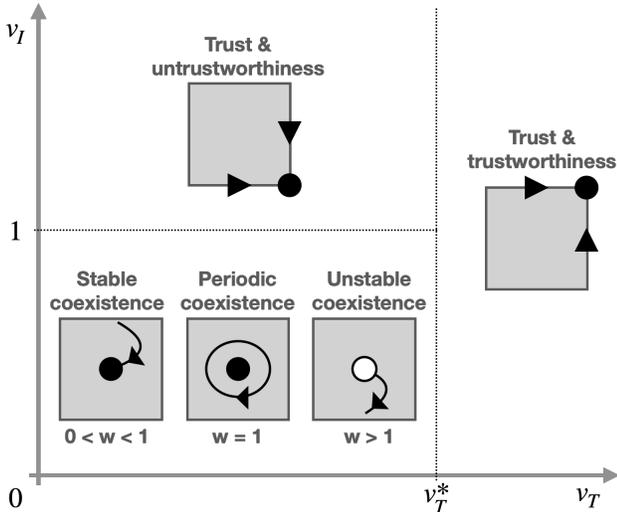


Fig. 3. Schematic summarising the evolutionary dynamics as a function of the incentive values  $v_I$  and  $v_T$ . On the boundaries of the state space, i.e. the unit square, we only show the stable equilibria and trajectories flowing into them. Non-generic cases (i.e.  $v_T = 0, v_T = v_T^*, v_I = 0$ , and  $v_I = 1$ ) are not shown.

after equilibration through the evolutionary dynamics (e.g. stable equilibria). Note  $\frac{\partial P}{\partial v_I} = -\frac{N_I(a-y_i)}{N_I+N_T} < 0$  since  $a > 1$  and  $y_i \leq 1$ . In other words, somewhat counterintuitively, the incentive given to investing investors,  $v_I$ , harms the overall social welfare in that the population average of the payoff decreases as  $v_I$  increases. Therefore, for any given  $(y_i, y_t)$ , one needs to minimise  $v_I$  to maximise  $P(y_i, y_t)$ .

1) *Optimal Payoff at (0, 0)*: The population average of the payoff at  $(0, 0)$  is given by

$$P(0, 0) = -\frac{aN_I v_I}{N_I + N_T}. \quad (22)$$

If  $(0, 0)$  is a stable equilibrium (i.e.  $0 \leq v_I < 1 \wedge v_T = 0$ ), then  $P(0, 0)$  is maximised at  $v_I = 0 \wedge v_T = 0$ .

2) *Optimal Payoff at (1, 0)*: The population average of the payoff at  $(1, 0)$  is

$$P(1, 0) = \frac{N_I(-av_I + v_I - 1)}{N_I + N_T} + \frac{N_T \left[ \frac{1-w^{N_I}}{N_T(1-w)} - av_T \right]}{N_I + N_T}. \quad (23)$$

We obtain  $\frac{\partial}{\partial v_T} P(1, 0) = -\frac{aN_T}{N_I + N_T} < 0$ . Therefore, if  $(1, 0)$  is an asymptotically stable equilibrium (i.e.  $v_I > 1 \wedge 0 \leq v_T <$

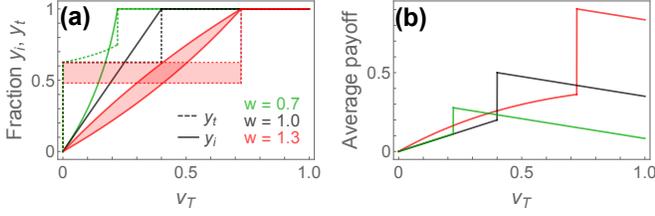


Fig. 4. Effects of the incentive to trustworthy trustees,  $v_T$ , and the nonlinearity in the payoff function,  $w$ , on the evolutionary outcomes in the NTG. We use the same parameter values as those used in Fig. 2 except for  $v_T$ . (a) Fractions of prosocial players as functions of the reward given to trustworthy trustees,  $v_T$ , in the equilibrium. We show the fraction of investing investors,  $y_i$ , and the fraction of trustworthy trustees,  $y_t$ . For  $w > 1$  and  $0 < v_T < v_T^*$ , the time averages of  $y_i$  and  $y_t$  do not converge. Therefore, we instead plot the ranges of asymptotic values of  $y_i$  and  $y_t$  by shaded regions. We observe that  $y_i$  increases as  $v_T$  increases when  $v_T < v_T^*$ . When  $v_T > v_T^*$ , the full trust  $y_i = 1$  and full trustworthiness  $y_t = 1$  evolve. (b) Population-averaged payoff,  $P$ , as a function of  $v_T$ . We observe that  $P$  increases as  $v_T$  increases when  $v_T < v_T^*$  and that  $P$  decreases as  $v_T$  increases when  $v_T > v_T^*$ . Note that the time average of  $P$  converges even if those of  $y_i$  and  $y_t$  do not. Panel (a) indicates that as  $w$  increases (i.e. from sub-linear to linear to super-linear), the evolution of trust and trustworthiness becomes more difficult. In other words, a higher value of  $v_T$  is necessary for attaining the same fraction of prosocial players when  $w$  is larger. In contrast, panel (b) indicates that the payoff of full trust and trustworthiness increases as  $w$  increases.

$v_T^*$ ), then  $P(1, 0)$  is maximised at  $v_I = 1 + \epsilon \wedge v_T = 0$ , where  $0 < \epsilon \ll 1$ .

3) *Optimal Payoff at (1, 1)*: The population average of the payoff at (1, 1) is

$$P(1, 1) = \frac{(1-a)(N_I v_I + N_T v_T)}{N_I + N_T} + \frac{2r(w^{N_I} - 1)}{(w-1)(N_I + N_T)}. \quad (24)$$

We obtain  $\frac{\partial}{\partial v_T} P(1, 1) = -\frac{(a-1)N_T}{N_I + N_T} < 0$ . If (1, 1) is an asymptotically stable equilibrium (i.e.  $v_T > v_T^*$ ), then  $P(1, 1)$  is maximised at  $v_I = 0 \wedge v_T = v_T^* + \epsilon$ .

4) *Optimal Payoff at Q or on Cycles around Q*: Recall that there exists a unique interior equilibrium  $\mathbf{Q}$  for  $0 \leq v_I < 1 \wedge 0 < v_T < v_T^*$ . For  $w < 1$ ,  $\mathbf{Q}$  is an asymptotically stable equilibrium and all the trajectories surrounding  $\mathbf{Q}$  converge to it. For  $w = 1$ , at which all the trajectories surrounding  $\mathbf{Q}$  form closed cycles, the time average of the population-mean payoff over the cycle is the same as the payoff at the equilibrium, i.e.,  $P(\mathbf{Q})$ ; see Appendix C for the proof. Therefore, seeking the optimal payoff at  $\mathbf{Q}$  is sufficient in both cases  $w < 1$  and  $w = 1$ . The population average of the payoff at  $\mathbf{Q}$  is given by

$$P(\mathbf{Q}) = \frac{N_T v_T - a(1-r)(N_I v_I + N_T v_T)}{(1-r)(N_I + N_T)}. \quad (25)$$

Note that  $P(\mathbf{Q})$  does not depend on  $w$ . We obtain  $\frac{\partial P(\mathbf{Q})}{\partial v_T} = \frac{N_T(1-a+ar)}{(1-r)(N_I + N_T)} > 0$  when  $r > r_0^* \equiv \frac{a-1}{a}$  and  $\frac{\partial P(\mathbf{Q})}{\partial v_T} < 0$  when  $r < r_0^*$ . Thus,  $P(\mathbf{Q})$  is monotonic as a function of  $v_T$  (Fig. 4b). For  $0 < w \leq 1$ , if  $\mathbf{Q}$  is asymptotically stable (i.e.,  $w < 1$ ) or neutrally stable (i.e.,  $w = 1$ ), then  $P(\mathbf{Q})$  is maximised at  $v_I = 0 \wedge v_T = v_T^* - \epsilon$  when  $r > r_0^*$  and at  $v_I = 0 \wedge v_T = 0 + \epsilon$  when  $r < r_0^*$ .

For  $w > 1$ , the time averages of  $y_i$  and  $y_t$  do not converge, but the time average of the payoff converges to

$$\bar{P}_{\text{hc}} = \frac{1}{(w-1) \left( \frac{(r+1)(N_I[1-r] + rN_T v_T)}{N_T v_T (r[w^{N_I} - 1] + N_I(w-1))} - \frac{r}{w^{N_I - 1}} \right) - a(N_I v_I + N_T v_T)} \frac{1}{N_I + N_T}, \quad (26)$$

where  $\bar{P}_{\text{hc}}$  is a convex combination of  $P(0, 0)$ ,  $P(0, 1)$ ,  $P(1, 0)$  and  $P(1, 1)$  as shown in Appendix D. Note that  $\frac{\partial \bar{P}_{\text{hc}}}{\partial v_I} = -\frac{aN_I}{N_I + N_T} < 0$  and that  $\bar{P}_{\text{hc}}$  is monotonic or has a local maximum as a function of  $v_T$ , as shown in Appendix E2. Therefore, given  $v_I = 0$ , the maximum of  $\bar{P}_{\text{hc}}(v_T)$  is either  $\bar{P}_{\text{hc}}(0 + \epsilon)$ ,  $\bar{P}_{\text{hc}}(v_T^{\text{hc}})$  or  $\bar{P}_{\text{hc}}(v_T^* - \epsilon)$ , where the local maximum of  $\bar{P}_{\text{hc}}(v_T)$  is at  $v_T = v_T^{\text{hc}} \equiv \frac{\left\{ \sqrt{aN_I(1-r^2)(w-1)[r(w^{N_I} - 1) + N_I(w-1)] - aN_I(1-r^2)(w-1)} \right\}}{aN_T r(w-1)(w^{N_I} - N_I w + N_I - 1)} \times (w^{N_I} - 1)$ .

5) *Comparison of the Optimal Payoff at the Different Equilibria*: We now compare the average payoff at the different equilibria. At each equilibrium, including the case of neutral and heteroclinic cycles, we denote by  $P^*$  the payoff maximised with respect to  $v_I$  and  $v_T$ . We compare  $P^*$  across the different equilibria to seek the overall maximum of the payoff and the associated optimal incentive.

For  $0 < w \leq 1$ , if  $r > r_1^* \equiv \frac{a-1}{a+1}$ , then the optimal payoff among the different equilibria is  $P^*(1, 1)$ ; if  $r < r_1^*$ , then the optimal payoff is  $P^*(0, 0)$ ; the associated optimal incentives are  $v_I = 0 \wedge v_T = v_T^* + \epsilon$  and  $v_I = 0 \wedge v_T = 0$ , respectively. For  $w > 1$ , as  $N_I \rightarrow \infty$  or  $w \rightarrow \infty$ , if  $r > r_2^* \equiv \frac{a}{a+1}$ , then the optimal payoff is  $P^*(1, 1)$ ; if  $r < r_2^*$ , then the optimal payoff is  $P^*(1, 0)$ ; the associated optimal incentives are  $v_I = 0 \wedge v_T = v_T^* + \epsilon$  and  $v_I = 1 + \epsilon \wedge v_T = 0$ , respectively. See Appendix E for the derivation of the optimal incentives. For relatively small values of  $w > 1$  and  $N_I \geq 2$ , the analytical derivation is not feasible and we instead numerically obtain the optimal incentives. Differently from the case of large  $N_I$  or  $w$ , the incentive yielding the heteroclinic cycle can realize the optimal payoff (Fig. 5). Note that copresence of incentives to investors and trustees (i.e.  $v_I > 0 \wedge v_T > 0$ ) is never optimal.

In summary, if the productivity of the prosocial strategies,  $r$ , is high enough relative to the fee rate  $a$ , the incentive leading to the full pro-sociality (i.e. full trust and full trustworthiness) is optimal. If the productivity is relatively low, the incentive leading to lower pro-sociality, including the case of the null incentive, is optimal.

#### F. Other Nonlinear Payoff Functions

To test the robustness of the results with respect to details of nonlinear payoff functions, we numerically examine evolutionary dynamics with nonlinear payoff functions that are different from but qualitatively similar to those given by Eq. (1). Specifically, we consider  $\log(k_i + 1)/\log(2)$  as a sub-linear payoff function that is qualitatively similar to Eq. (1) with  $0 < w < 1$  and  $\exp(0.7k_i) - 1$  as a super-linear payoff

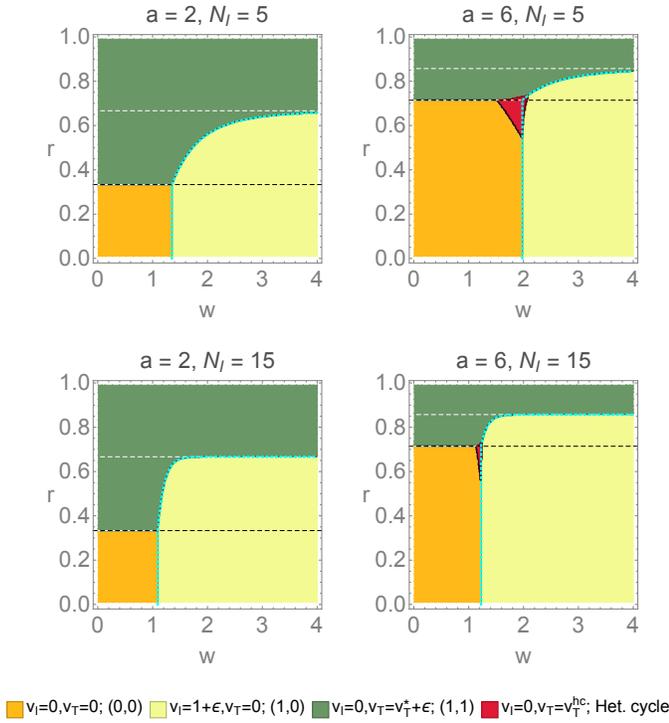


Fig. 5. Optimal incentives and associated evolutionary outcomes. Each colored region shows the parameter region in which the associated stable equilibrium or the heteroclinic cycle yields the largest population average of the payoff given by Eq. (21). Among the two horizontal dashed lines, the lower and upper ones indicate  $r = r_1^* = \frac{a-1}{a+1}$  and  $r = r_2^* = \frac{a}{a+1}$ , respectively. The dotted curve indicates  $r(w) = r_2^* - \frac{aN_I(w-1)}{(a+1)(w^{N_I}-1)}$ , where  $P^*(1, 1) = P^*(1, 0)$ ; note that  $r(w) \rightarrow r_2^*$  as  $w$  increases. The vertical dotted line indicates  $w = w^* > 1$  that we obtained by numerically solving  $P^*(0, 0) = P^*(1, 0)$ . As the fee rate  $a$  increases, the parameter region in which  $(0, 0)$  is optimal with the null incentive (in orange) and the region in which the heteroclinic cycle is optimal with a positive incentive (in red) become larger. As  $N_I$  or  $w$  increases, the border between parameter region in which  $P^*(1, 1)$  is optimal (in dark green) and that in which  $P^*(1, 0)$  is optimal (in light yellow) converges to  $r = r_2^*$ , which we have analytically derived in the limit  $N_I \rightarrow \infty$  or  $w \rightarrow \infty$ . For a larger fee rate,  $a$ , or a larger size of the investor group,  $N_I$ , the incentive yielding full trust and trustworthiness is optimal for a smaller parameter region (i.e. the green regions in the figure).

function that is qualitatively similar to Eq.(1) with  $w > 1$ . Figure 6 indicates that each of these payoff functions yields qualitatively the same evolutionary dynamics as those obtained with Eq. (1).

#### IV. DISCUSSION

The  $N$ -player generalisation of a TG game proposed in Ref. [27] assumes that an investor always invests. Therefore, their NTG is structurally different from both the two-player TG and our NTG. It may be instead called the trustworthiness game in that the payoff of the game is entirely determined by the strategy of a trustee. The ultimatum game (UG) and the dictator game (DG) already have a parallel to this distinction between the TG and the trustworthiness game. The UG involves a non-simultaneous interaction on resource split between a proposer and a responder [35]. The simplest variant of the UG assumes two options for each role: for a proposer

to propose an unfair split in favour of the proposer or a fair split, and for a responder to accept or reject the proposal. If the responder accepts, both the proposer and responder obtain the proposed payoffs. If the response rejects, both players get nothing. The DG is similar to the UG except that a responder has no option other than to accept any proposal made by the proposer. Hence, the payoff entirely depends on what a proposer does and thus the proposer is called a dictator. The DG is related to but structurally different from the UG, and therefore the DG has been analysed on its own [36][37][38]. In the UG, the reputation mechanism, which is equivalent to a responder refusing an unfair split, can lead a proposer to offer a fair one [35]. However, the reputation mechanism cannot work for the DG since a responder has no option of refusing any split. Our NTG is of the UG type in that it allows the investor an option not to invest, which has enabled us to investigate the evolution of trust as well as trustworthiness.

The evolutionary game dynamics in Ref. [27] assumes role-unaware imitation, which allows imitation between the different roles and leads to the cease of game playing. The justification of this assumption is unclear. To the best of our knowledge, this type of game dynamics has not been used prior to Ref. [27], regardless of two-player or  $N$ -player games. In fact, there have been two canonical approaches to modelling evolutionary dynamics of non-simultaneous games. One approach is to assume that each player plays each role half of the time and imitates others in a role-aware manner [35][39][40]. The player's strategy is then a tuple consisting of the strategies under the different roles (e.g. one as an investor and the other as a trustee). This symmetrisation probably better characterises scenarios in which each player has multiple roles, and thus, the payoff of the player is the average of the payoffs from the different roles. For instance, a bank can lend money to or borrow money from other banks, playing two roles, as a lender/investor and a borrower/trustee. By this symmetrisation, one can consider the TG using a single population and the corresponding replicator dynamics [30][40]. The same approach has also been used for other asymmetric games such as the UG [35][40]. Developing and analyzing NTGs with this symmetrisation method is an open question. A second approach is to fix the two roles such that players can imitate others in their own role only [19][41]. We took this approach to formulate an asymmetric NTG, which is a faithful generalisation of a previously proposed two-player TG [19]. Then, differently from the previous work allowing the imitation between the different roles and hence leading to the extinction of investors [27], we found that investors do not perish but evolve not to trust trustees unless an incentive is in place.

The payoff in Ref. [27] is a linear function of the number of investing investors, which is also inherited in its follow-up studies [31][32][33]. With a linear payoff function, any  $N$ -player game is equivalent to a sum of two-player games and thus the evolutionary outcome of the former is similar to that of the latter. In  $N$ -player games, however, unlike two-player games, nonlinear payoff functions can yield evolutionary outcomes that are qualitatively different from those of linear ones. We have introduced nonlinear payoff functions

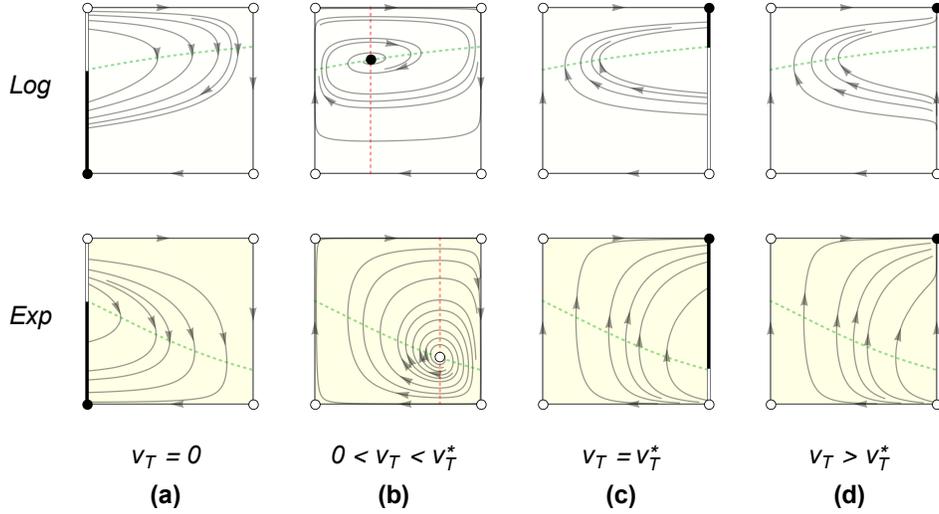


Fig. 6. Robustness of the evolutionary dynamics with respect to details of nonlinear payoff functions. The parameters are the same as those used in Fig. 2. The top panels show that a sub-linear payoff function  $\log(k_i + 1) / \log(2)$  leads to evolutionary dynamics similar to that for Eq. (1) with  $w = 0.7$ , which is presented in the top panels in Fig. 2. The bottom panels show that a super-linear payoff function  $\exp(0.7k_i) - 1$  leads to evolutionary dynamics similar to that for Eq. (1) with  $w = 1.3$ , which is presented in the bottom panels in Fig. 2.

in the asymmetric NTG. Even with the nonlinear payoff functions, we have found that it is more challenging for pro-social behaviours to evolve in the asymmetric NTG than in the PGG. The PGG is one of the most widely studied  $N$ -player games [42]. With linear payoff functions, the PGG becomes a dominance game for which anti-social behaviour (i.e. defection) dominates pro-social behaviour in terms of the payoff value and hence only anti-social behaviour evolves. With the non-linear payoff functions of the same form used in the present paper, the PGG becomes either a coexistence game or a coordination game for which prosocial behaviour can evolve [4]. Therefore, incentives have been applied only to the linear PGGs but not the nonlinear PGGs; see Ref. [43] for a review. In the asymmetric NTG with fixed roles, however, we have found that the nonlinear payoff functions are not sufficient for pro-social behaviour to evolve and an additional mechanism such as an incentive is required. We have found that the incentive to trustworthy trustees can be sufficient for the full pro-sociality to evolve in both investor and trustee populations, i.e., the full trust (i.e. investment) and the full trustworthiness. An intuitive explanation of this result is as follows. If the fraction of trustworthy trustees is high enough, the payoff of investing investors is higher than that of non-investing ones and thus investing investors evolve. Hence, if the incentive to trustworthy trustees is large enough for them to evolve, then it also yields the evolution of investing investors.

With the nonlinear payoff function given by Eq. (1), one can express the discount (i.e., sub-linear) and synergy (i.e., super-linear) effects by tuning the single parameter  $w$ . This payoff function is advantageous because it allows us to analytically examine the evolutionary dynamics for arbitrary group sizes  $N_I$  and  $N_T$ . However, our results are not confined to this particular form of payoff function. We ran numerical simulations with different payoff functions to support that our results are robust with respect to details of the nonlinearity

of the payoff function. We remark that, unlike with Eq. (1), different nonlinear payoff functions require separate analyses of evolutionary dynamics for each combination of the values of  $N_I$  and  $N_T$  in general. Specifically, one needs to numerically find the interior equilibrium and carry out the linear stability analysis for each given  $N_I$  and  $N_T$ .

Given an investing investor, the two-player TG creates a social dilemma [19][27]. The total wealth (i.e. the sum of the payoffs of an investing investor and a trustee) depends on the strategy of a trustee. Although a self-interested (i.e. untrustworthy) trustee earns higher than a pro-social (i.e. trustworthy) trustee does, the former leads to a lower total wealth ( $= 0$ ) than the latter does ( $= 2r$ ) (Fig. 1a). Our NTG preserves the nature of a social dilemma. For a linear payoff function, given the number of investing investors,  $k_i$ , if all trustees in a group are self-interested, they earn more than any pro-social trustees would. However, the former leads to a lower total wealth ( $= 0$ ) than the latter does ( $= 2rk_i$ ).

Most previous studies on institutional incentives have focused on which incentives promote prosocial behaviours the best [44][45][33]. However, a better criterion for the success of an incentive may be the population average of payoff at the evolutionarily stable state [46]. Thus, we have sought the optimal incentive that yields the highest payoff, taking into consideration the operating cost of managing incentives. We have found that the incentive leading to the most prosocial behaviour (i.e. full trust and full trustworthiness) often yields the highest payoff but not always. When the productivity of the prosocial behaviours is not high enough, the incentive leading to less prosocial behaviours (e.g. combination of full trust and null trustworthiness) can yield the highest payoff; even the null incentive leading to null trust and null trustworthiness can be optimal when the operating cost of managing incentives outweighs benefits from prosocial behaviours. A limitation of our incentive scheme is to have assumed that an incentive

is tailored to individual players while the game is played in groups. Although this type of the individually targeted incentive is widely used for  $N$ -player games [6][45][47][48], it may be less feasible than it is for two-player games, in which actions of the individual players are more easily identified than in  $N$ -player games. Relaxing this assumption is worthwhile investigation. For instance, a diluted incentive scheme, which provides an incentive to a group, may be more feasible for  $N$ -player games. In such an incentive scheme, all individuals in a group receive the same incentive by construction, and whether a group receives an incentive is determined based on aggregated information such as the proportion of trustworthy trustees in the group.

In summary, we started by noting that the  $N$ -player TG in Ref. [27] is structurally different from the TG and proposed an asymmetric  $N$ -player TG with two fixed roles. With this setup, it is more challenging for pro-social strategies to evolve than in the celebrated PGG. Nonetheless, we showed that incentives provided to trustees can cost-effectively promote the evolution of trust and trustworthiness among self-interested players. We also showed that nonlinear payoff functions in the  $N$ -player TG yield a richer set of evolutionary dynamics and the associated optimal incentives than linear payoff functions. We hope that our contribution paves the way for further studies of  $N$ -player TGs and their variations such as the symmetrisation of asymmetric  $N$ -player TGs, the impacts of structured populations [49][50], repeated interactions on the evolution of trust/trustworthiness, and stochastic evolutionary dynamics in finite populations. There can be different generalisations of the two-player NTG each of which recovers the two-player TG when  $N_I = N_T = 1$ ; such generalisations are interesting to explore. Applications of  $N$ -player TGs are also worthwhile seeking; for instance, multi-hop relay in wireless sensors or ad hoc networks could be mapped to an  $N$ -player TG among self-interested nodes [51][52].

## APPENDIX

### A. Derivation of Eq. (9)

We obtain

$$\begin{aligned}
P_i^o &= \sum_{m_i=0}^{N_I-1} \binom{N_I-1}{m_i} y_i^{m_i} (1-y_i)^{N_I-1-m_i} \\
&\quad \times \sum_{m_t=0}^{N_T} \binom{N_T}{m_t} y_t^{m_t} (1-y_t)^{N_T-m_t} \Pi_i^o(m_i+1, m_t) \\
&= \sum_{m_i=0}^{N_I-1} \binom{N_I-1}{m_i} y_i^{m_i} (1-y_i)^{N_I-1-m_i} \\
&\quad \times \sum_{m_t=0}^{N_T} \binom{N_T}{m_t} y_t^{m_t} (1-y_t)^{N_T-m_t} \left[ \frac{m_t}{N_T} \frac{r(1-w^{m_i+1})}{(m_i+1)(1-w)} \right. \\
&\quad \quad \left. + \left(1 - \frac{m_t}{N_T}\right) \cdot (-1) \right] \\
&= \sum_{m_i=0}^{N_I-1} \binom{N_I-1}{m_i} y_i^{m_i} (1-y_i)^{N_I-1-m_i} \\
&\quad \times \left[ \frac{1}{m_i+1} r \frac{1-w^{m_i+1}}{1-w} + 1 \right] y_t - 1
\end{aligned}$$

$$\begin{aligned}
&= \frac{r y_t}{N_I(1-w)} \left[ \sum_{m_i=0}^{N_I-1} \binom{N_I}{m_i+1} y_i^{m_i} (1-y_i)^{N_I-1-m_i} \right. \\
&\quad \left. \times (1-w^{m_i+1}) \right] + y_t - 1 \\
&= \frac{r}{N_I(1-w)} \frac{y_t}{y_i} \left[ 1 - (1+(w-1)y_i)^{N_I} \right] + y_t - 1,
\end{aligned} \tag{A.27}$$

where we have assumed that  $y_i \neq 0$  and used the expression of the mean of a binomial distribution  $\sum_{m_t=0}^{N_T} \binom{N_T}{m_t} y_t^{m_t} (1-y_t)^{N_T-m_t} m_t = N_T y_t$  and the relationship  $\binom{N_I-1}{m_i} \frac{1}{m_i+1} = \frac{1}{N_I} \binom{N_I}{m_i+1}$ . To show the last equality in Eq. (A.27), with substitution  $k_i \equiv m_i + 1$ , we used

$$\begin{aligned}
&\sum_{m_i=0}^{N_I-1} \binom{N_I}{m_i+1} y_i^{m_i} (1-y_i)^{N_I-1-m_i} (1-w^{m_i+1}) \\
&= \sum_{k_i=1}^{N_I} \binom{N_I}{k_i} y_i^{k_i-1} (1-y_i)^{N_I-k_i} (1-w^{k_i}) \\
&= \frac{1}{y_i} \sum_{k_i=1}^{N_I} \binom{N_I}{k_i} y_i^{k_i} (1-y_i)^{N_I-k_i} (1-w^{k_i}) \\
&= \frac{1}{y_i} \left[ \sum_{k_i=0}^{N_I} \binom{N_I}{k_i} y_i^{k_i} (1-y_i)^{N_I-k_i} (1-w^{k_i}) \right. \\
&\quad \left. - \binom{N_I}{0} y_i^0 (1-y_i)^{N_I} (1-w^0) \right] \\
&= \frac{1}{y_i} \left[ \sum_{k_i=0}^{N_I} \binom{N_I}{k_i} y_i^{k_i} (1-y_i)^{N_I-k_i} (1-w^{k_i}) \right] \\
&= \frac{1}{y_i} \left[ (y_i+1-y_i)^{N_I} - \sum_{k_i=0}^{N_I} \binom{N_I}{k_i} y_i^{k_i} (1-y_i)^{N_I-k_i} w^{k_i} \right] \\
&= \frac{1}{y_i} \left[ 1 - \sum_{k_i=0}^{N_I} \binom{N_I}{k_i} (w y_i)^{k_i} (1-y_i)^{N_I-k_i} \right] \\
&= \frac{1}{y_i} \left[ 1 - (w y_i + 1 - y_i)^{N_I} \right] \\
&= \frac{1}{y_i} \left[ 1 - (1+(w-1)y_i)^{N_I} \right].
\end{aligned} \tag{A.28}$$

We have  $P_i = P_i^o + v_I - a v_I$ , where  $P_i^o$  is given by Eq. (A.27). We can similarly derive  $P_t$  and  $P_u$ .

### B. Existence and Stability of the Equilibria

One can deduce the signs of the two eigenvalues  $\lambda_1$  and  $\lambda_2$  of the Jacobian matrix,  $J$ , at an equilibrium by its determinant and trace, which are equal to  $\lambda_1 \lambda_2$  and  $\lambda_1 + \lambda_2$ , respectively. We denote by  $\text{Det}|_{\mathbf{y}}$  and  $\text{Tr}|_{\mathbf{y}}$  the determinant and trace, respectively, of  $J$  evaluated at  $\mathbf{y} \in [0, 1]^2$ . Especially, the asymptotical stability of an equilibrium requires  $\lambda_1 < 0$  and  $\lambda_2 < 0$ , which lead to  $\text{Det}|_{\mathbf{y}} > 0$  and  $\text{Tr}|_{\mathbf{y}} < 0$ . We determine the stability of each equilibrium as follows.

1) (0, 0): The Jacobian matrix at  $(y_i, y_t) = (0, 0)$  is given by

$$J_{(0,0)} = \begin{pmatrix} v_I - 1 & 0 \\ 0 & v_T \end{pmatrix}. \tag{A.29}$$

We obtain

$$\text{Det}|_{(0,0)} = (v_I - 1)v_T \tag{A.30}$$

and

$$\text{Tr}|_{(0,0)} = v_I + v_T - 1. \quad (\text{A.31})$$

If  $0 \leq v_I < 1 \wedge v_T = 0$ , then  $\text{Det}|_{(0,0)} = 0 \wedge \text{Tr}|_{(0,0)} < 0$  such that  $(0, 0)$  is stable but not asymptotically stable. Otherwise,  $(0, 0)$  is unstable.

2)  $(0, 1)$ : The Jacobian at  $(0, 1)$  is given by

$$J_{(0,1)} = \begin{pmatrix} r + v_I & 0 \\ 0 & -v_T \end{pmatrix}. \quad (\text{A.32})$$

We obtain

$$\text{Det}|_{(0,1)} = -v_T(r + v_I) \quad (\text{A.33})$$

and

$$\text{Tr}|_{(0,1)} = r + v_I - v_T. \quad (\text{A.34})$$

If  $v_T = 0$ , then  $\text{Det}|_{(0,1)} = 0 \wedge \text{Tr}|_{(0,1)} > 0$  such that  $(0, 1)$  is unstable. If  $v_T > 0$ , then  $\text{Det}|_{(0,1)} < 0$  such that  $(0, 1)$  is unstable.

3)  $(1, 0)$ : The Jacobian at  $(1, 0)$  is given by

$$\begin{aligned} J_{(1,0)} &= \begin{pmatrix} 1 - v_I & 0 \\ 0 & v_T - \frac{(1-r)(w^{N_I}-1)}{N_T(w-1)} \end{pmatrix} \\ &= \begin{pmatrix} 1 - v_I & 0 \\ 0 & v_T - v_T^* \end{pmatrix}. \end{aligned} \quad (\text{A.35})$$

We obtain

$$\text{Det}|_{(1,0)} = (v_T - v_T^*)(1 - v_I) \quad (\text{A.36})$$

and

$$\text{Tr}|_{(1,0)} = v_T - v_T^* + 1 - v_I, \quad (\text{A.37})$$

where

$$v_T^* = \frac{(1-r)(w^{N_I}-1)}{N_T(w-1)} > 0. \quad (\text{A.38})$$

If  $v_I > 1 \wedge v_T < v_T^*$ , then  $\text{Det}|_{(1,0)} > 0 \wedge \text{Tr}|_{(1,0)} < 0$  such that  $(1, 0)$  is asymptotically stable. If  $(v_I > 1 \wedge v_T = v_T^*) \vee (v_I = 1 \wedge v_T < v_T^*)$ , then  $\text{Det}|_{(1,0)} = 0 \wedge \text{Tr}|_{(1,0)} < 0$  such that  $(1, 0)$  is stable but not asymptotically stable. Otherwise,  $(1, 0)$  is unstable.

4)  $(1, 1)$ : The Jacobian at  $(1, 1)$  is given by

$$\begin{aligned} J_{(1,1)} &= \begin{pmatrix} -\frac{r(w^{N_I}-1)}{N_I(w-1)} - v_I & 0 \\ 0 & \frac{(1-r)(w^{N_I}-1)}{N_T(w-1)} - v_T \end{pmatrix} \\ &= \begin{pmatrix} -\frac{N_T r v_T^*}{N_I(1-r)} - v_I & 0 \\ 0 & v_T^* - v_T \end{pmatrix}. \end{aligned} \quad (\text{A.39})$$

We obtain

$$\text{Det}|_{(1,1)} = (v_T - v_T^*) \left[ v_I + \frac{N_T r v_T^*}{N_I(1-r)} \right] \quad (\text{A.40})$$

and

$$\text{Tr}|_{(1,1)} = \left[ 1 - \frac{N_T r}{N_I(1-r)} \right] v_T^* - v_I - v_T. \quad (\text{A.41})$$

We obtain  $\text{sign}(\text{Det}|_{(1,1)}) = \text{sign}(v_T - v_T^*)$  since  $v_I + \frac{N_T r v_T^*}{N_I(1-r)} > 0$ , which is guaranteed by  $0 < r < 1$ ,  $v_I \geq 0$  and  $v_T^* > 0$ . If  $v_T > \left[ 1 - \frac{N_T r}{N_I(1-r)} \right] v_T^* - v_I$ , then  $\text{Tr}|_{(1,1)} < 0$ . We also note  $\left[ 1 - \frac{N_T r}{N_I(1-r)} \right] v_T^* - v_I \leq \left[ 1 - \frac{N_T r}{N_I(1-r)} \right] v_T^* < v_T^*$ .

Therefore, if  $v_T > v_T^*$ , then  $\text{Det}|_{(1,1)} > 0 \wedge \text{Tr}|_{(1,1)} < 0$  such that  $(1, 1)$  is asymptotically stable.

If  $v_T = v_T^*$ , then  $\text{Det}|_{(1,1)} = 0 \wedge \text{Tr}|_{(1,1)} < 0$ . Therefore,  $(1, 1)$  is stable but not asymptotically stable.

If  $v_T < v_T^*$ , then  $\text{Det}|_{(1,1)} < 0$ . Therefore,  $(1, 1)$  is unstable.

5) *Interior equilibrium Q*: We show that there exists a unique interior equilibrium **Q** if and only if  $0 < v_T < v_T^* = \frac{(1-r)(w^{N_I}-1)}{N_T(w-1)}$  and  $v_I < 1$ . The internal equilibrium, if it exists, is located at the intersection of the nullclines  $P_t(y_i, y_t) - P_u(y_i, y_t) = 0$  and  $P_i(y_i, y_t) - P_n(y_i, y_t) = 0$  with  $0 < y_i < 1 \wedge 0 < y_t < 1$ . Let us investigate the two nullclines one by one.

Because  $P_t - P_u = \frac{(1-r)\{1-[(w-1)y_i+1]^{N_I}\}}{N_T(w-1)} + v_T$  does not depend on  $y_t$ , the nullcline  $P_t - P_u = 0$  is of the form  $y_i = \text{constant}$ . Specifically,  $P_t - P_u = 0$  leads to  $y_i = y_{i,Q} \equiv \frac{d^{1/N_I}-1}{w-1}$ , where  $d = 1 + \frac{N_T v_T (w-1)}{1-r}$ . We obtain  $\frac{d}{dy_i} [P_t - P_u] = -\frac{N_I(1-r)\{(w-1)y_i+1\}^{N_I-1}}{N_T} < 0$ . Therefore, if and only if  $0 < v_T < v_T^*$ , then  $P_t(0, y_t) - P_u(0, y_t) = v_T > 0$  and  $P_t(1, y_t) - P_u(1, y_t) = v_T - v_T^* < 0$  such that the nullcline  $P_t - P_u = 0$  (i.e.  $y_i = y_{i,Q}$ ) exists with  $0 < y_{i,Q} < 1$ .

To examine the other nullcline, we look into  $P_i - P_n = \frac{r y_i \{1-[1+(w-1)y_i]^{N_I}\}}{N_I(1-w)y_i} + y_t - 1 + v_I$ . In fact,  $\frac{\partial}{\partial y_t} [P_i - P_n] > 0$  and  $P_i(y_i, 1) - P_n(y_i, 1) > 0$  hold true for  $0 < y_i < 1$ , which we will show later. Therefore, if and only if  $v_I < 1$ , then  $P_i(y_i, 0) - P_n(y_i, 0) = v_I - 1 < 0$  such that the nullcline  $P_i - P_n = 0$  exists in the range  $0 < y_t < 1$ . Note that the nullcline  $P_i - P_n = 0$  can be represented by  $y_t = g(y_i)$  because there exists a unique  $y_t$  satisfying  $P_i - P_n = 0$  for any  $y_i$ .

We now show  $\frac{\partial}{\partial y_t} [P_i - P_n] = \frac{r\{1-[1+(w-1)y_i]^{N_I}\}}{N_I(1-w)y_i} + 1 > 0$  for  $0 < y_i < 1$ . If  $0 < w < 1$ , then we obtain  $0 < 1 + (w-1)y_i < 1$  such that  $1 - [1 + (w-1)y_i]^{N_I}$  and  $1 - w$  are both positive. If  $w > 1$ , then  $1 < 1 + (w-1)y_i$  such that  $1 - [1 + (w-1)y_i]^{N_I}$  and  $1 - w$  are both negative. If  $w = 1$ , then  $\frac{\partial}{\partial y_t} [P_i - P_n] = \frac{\lim_{w \rightarrow 1} r\{1-[1+(w-1)y_i]^{N_I}\}}{\lim_{w \rightarrow 1} N_I(1-w)y_i} + 1 = r + 1 > 0$ . Therefore, we have proved  $\frac{\partial}{\partial y_t} [P_i - P_n] > 0$  for any  $w$ .

We now show  $P_i(y_i, 1) - P_n(y_i, 1) > 0$ . If  $w \neq 1$ , then we obtain  $P_i(y_i, 1) - P_n(y_i, 1) = \frac{r\{1-[1+(w-1)y_i]^{N_I}\}}{N_I(1-w)y_i} + v_I > 0$ . If  $w = 1$ , then we obtain  $P_i(y_i, 1) - P_n(y_i, 1) = \frac{\lim_{w \rightarrow 1} r\{1-[1+(w-1)y_i]^{N_I}\}}{\lim_{w \rightarrow 1} N_I(1-w)y_i} + v_I = r + v_I > 0$ . Therefore,  $P_i(y_i, 1) - P_n(y_i, 1) > 0$  holds true for any  $w$ .

Finally, these results imply that there is a unique intersection of  $y_i = y_{i,Q}$  and  $y_t = g(y_i)$  satisfying  $0 < y_i < 1 \wedge 0 < y_t < 1$ , which is an interior equilibrium **Q**.

We now analyse the stability of the interior equilibrium **Q**. The Jacobian at **Q** is given by

$$J_Q = \begin{pmatrix} J_{11}^Q & J_{12}^Q \\ J_{21}^Q & 0 \end{pmatrix}, \quad (\text{A.42})$$

where  $J_{11}^Q = \frac{N_I w d - d^{1/N_I} \{d[(N_I-1)w+N_I]+w\} + [d(N_I-1)+1]d^{2/N_I}}{(w-1)\{d^{1/N_I}[N_I-(d-1)r]-N_I d^{2/N_I}\}}$ ,  $r(1 - v_I)$ ,  $J_{12}^Q = -\frac{\{N_I d^{1/N_I} + [(d-1)r - N_I]\}(d^{1/N_I} - w)}{N_I(w-1)^2}$ , and  $J_{21}^Q = -\frac{d^{1-1/N_I}(d^{1/N_I}-1)\{N_I v_I d^{1/N_I} + [(d-1)r - N_I v_I]\}}{N_T\{N_I d^{1/N_I} + [(d-1)r - N_I]\}^2} \times$

$N_I^2(r-1)(v_I-1)$ . We first show  $\text{Det}|_Q > 0$  and  $\text{sign}(\text{Tr}|_Q) = \text{sign}(w-1)$ .

For  $w \neq 1$ , we have  $\text{Det}|_Q = (v_I-1)N_I(1-r)d^{1-1/N_I} \frac{(d^{1/N_I}-1)(d^{1/N_I}-w)[N_I v_I(d^{1/N_I}-1)+(d-1)r]}{N_T(w-1)^2[N_I(d^{1/N_I}-1)+(d-1)r]}$ .

We note that  $\frac{N_I v_I(d^{1/N_I}-1)+(d-1)r}{N_I(d^{1/N_I}-1)+(d-1)r}$  is positive because  $\text{sign}(d-1) = \text{sign}(d^{1/N_I}-1)$ . Since  $0 < v_T < v_T^* = \frac{(1-r)(w^{N_I}-1)}{N_T(w-1)}$  for the existence of the interior equilibrium, we have  $d = 1 + \frac{N_T v_T(w-1)}{1-r} = 1 + s(w^{N_I}-1)$ , where  $v_T = s v_T^*$  and  $0 < s < 1$ . Therefore, we obtain  $w^{N_I} - d = (1-s)(w^{N_I}-1) \implies \text{sign}(w^{N_I}-d) = \text{sign}(w^{N_I}-1) \implies \text{sign}(w-d^{1/N_I}) = \text{sign}(w-1) \implies (w < d^{1/N_I} < 1) \vee (1 < d^{1/N_I} < w) \implies (d^{1/N_I}-1)(d^{1/N_I}-w) < 0 \implies \text{sign}((d^{1/N_I}-1)(d^{1/N_I}-w)) = -1$ . For  $d \neq 1$ , we obtain  $\text{sign}(\text{Det}|_Q) = \text{sign}(1-v_I)$  because  $\text{sign}(\text{Det}|_Q) = \text{sign}(v_I-1) \text{sign}[(d^{1/N_I}-1)(d^{1/N_I}-w)] \times \text{sign}\left(\frac{N_I v_I(d^{1/N_I}-1)+(d-1)r}{N_I(d^{1/N_I}-1)+(d-1)r}\right) = (-1) \cdot (-1) \cdot 1 = 1$ .

Recall that  $0 \leq v_I < 1$  is required for the existence of  $\mathbf{Q}$ . For  $w = 1$ , we obtain  $\text{Det}|_Q = \frac{(1-v_I)v_T(r+v_I)[N_I(1-r)-N_T v_T]}{N_I(1-r)^2} > 0$  since  $0 < v_T < v_T^* = \frac{\lim_{w \rightarrow 1}(1-r)(w^{N_I}-1)}{\lim_{w \rightarrow 1} N_T(w-1)} = \frac{N_I(1-r)}{N_T}$  is required for the existence of  $\mathbf{Q}$ . Hence, we have shown  $\text{Det}|_Q > 0$  or  $\text{sign}(\text{Det}|_Q) = 1$  regardless of the  $w$  value.

For  $w \neq 1$ , we obtain  $\text{Tr}|_Q = \frac{r(1-v_I)d^{-1/N_I} \{d^{1/N_I}[(d(N_I-1)+1]-dN_I]\frac{w-d^{1/N_I}}{w-1}}{N_I(d^{1/N_I}-1)+(d-1)r}$ . We obtain  $\frac{\text{sign}(w-d^{1/N_I})}{\text{sign}(w-1)} = 1$  since  $\text{sign}(w-d^{1/N_I}) = \text{sign}(w-1)$  as already shown. For  $d \neq 1$ , we have  $q(d) \equiv d^{1/N_I} [d(N_I-1)+1] - dN_I > 0$  since  $q(1) = 0$  is the global minimum of  $q(d)$  for  $d > 0$ , the latter of which can be shown as follows. First,  $q(1)$  is a local minimum of  $q(d)$  since  $\left. \frac{\partial q}{\partial d} \right|_{d=1} = \left\{ \frac{d^{1/N_I-1} [d(N_I^2-1)+1]}{N_I} - N_I \right\} \Big|_{d=1} = 0$  and  $\left. \frac{\partial^2 q}{\partial d^2} \right|_{d=1} = \frac{(N_I-1)d^{1/N_I-2}(dN_I+d-1)}{N_I^2} \Big|_{d=1} = \frac{N_I-1}{N_I} > 0$ . Second,  $d = d^* \equiv \frac{1}{N_I+1} \in (0, 1)$  is the only inflection point of the function  $q(d)$  for  $d > 0$  since  $\frac{\partial^2 q}{\partial d^2} < 0$  for  $d < d^*$ ,  $\frac{\partial^2 q}{\partial d^2} = 0$  at  $d = d^*$ , and  $\frac{\partial^2 q}{\partial d^2} > 0$  for  $d > d^*$ . Therefore, there is no local minimum in  $d \leq d^*$  and at most one local minimum in  $d > d^*$ , which is at  $d = 1$ . Third, we obtain  $q(0) = 0 = q(1)$ . Hence,  $q(1) = 0$  is the global minimum of  $q(d)$  for  $d > 0$ . We obtain  $\text{sign}(N_I(d^{1/N_I}-1) + (d-1)r) = \text{sign}(w-1)$  since  $\text{sign}(d^{1/N_I}-1) = \text{sign}(d-1) = \text{sign}(w-1)$ . It follows that  $\text{sign}(\text{Tr}|_Q) = \text{sign}(1-v_I) \frac{\text{sign}(d^{1/N_I}[(d(N_I-1)+1]-dN_I]) \text{sign}(w-d^{1/N_I})}{\text{sign}(N_I(d^{1/N_I}-1)+(d-1)r) \text{sign}(w-1)} = 1 \cdot \frac{1}{\text{sign}(w-1)} \cdot 1 = \text{sign}(w-1)$ . For  $w = 1$ , it holds true that  $\text{sign}(\text{Tr}|_Q) = 0$  because  $\text{Tr}|_Q = 0$ . Hence, we have shown  $\text{sign}(\text{Tr}|_Q) = \text{sign}(w-1)$  regardless of the  $w$  value.

For  $w < 1$ , we obtain  $\text{Det}|_Q > 0 \wedge \text{Tr}|_Q < 0$  such that  $\mathbf{Q}$  is asymptotically stable. For  $w > 1$ , we obtain  $\text{Det}|_Q > 0 \wedge \text{Tr}|_Q > 0$  such that  $\mathbf{Q}$  is unstable. For  $w = 1$ , we obtain  $\text{Det}|_Q > 0$  and  $\text{Tr}|_Q = 0$ . In this case, the discriminant  $D = (\text{Tr}|_Q)^2 - 4\text{Det}|_Q < 0$  and  $\text{Tr}|_Q = 0$ , which implies that the eigenvalues are purely imaginary. Therefore,  $\mathbf{Q}$  is neutrally stable and the trajectories cycle around it.

6)  $y_i = 0$ : We find that  $(0, y_t)$ , where  $0 < y_t < 1$ , is a line of equilibria if and only if  $v_T = 0$ . In this case, the Jacobian at  $(0, y_t)$  is given by

$$J_{(0, y_t)} = \begin{pmatrix} r y_t + v_I + y_t - 1 & 0 \\ -\frac{N_I(r-1)(y_t-1)y_t}{N_T} & 0 \end{pmatrix}. \quad (\text{A.43})$$

We obtain  $\text{Det}|_{(0, y_t)} = 0$  and  $\text{Tr}|_{(0, y_t)} = (r+1)y_t + v_I - 1$ . If  $v_T = 0 \wedge y_t < \frac{1-v_I}{r+1}$ , then  $\text{Tr}|_{(0, y_t)} < 0$  such that  $(0, y_t)$  is stable but not asymptotically stable. If  $v_T = 0 \wedge y_t > \frac{1-v_I}{r+1}$ , then  $\text{Tr}|_{(0, y_t)} > 0$  such that  $(0, y_t)$  is unstable.

7)  $y_i = 1$ : We find that  $(1, y_t)$ , where  $0 < y_t < 1$ , is a line of equilibria if and only if  $v_T = v_T^*$ . In this case, the Jacobian at  $(1, y_t)$  is given by

$$J_{(1, y_t)} = \begin{pmatrix} -\frac{r y_t (w^{N_I}-1)}{N_I(w-1)} - v_I - y_t + 1 & 0 \\ -\frac{N_I(r-1)(y_t-1)y_t w^{N_I-1}}{N_T} & 0 \end{pmatrix}. \quad (\text{A.44})$$

We obtain  $\text{Det}|_{(1, y_t)} = 0$  and  $\text{Tr}|_{(1, y_t)} = -\frac{r y_t (w^{N_I}-1)}{N_I(w-1)} - v_I - y_t + 1$ . If  $v_T = v_T^* \wedge 0 \leq v_I < 1 \wedge y_t > \frac{N_I(1-v_I)(w-1)}{r(w^{N_I}-1)+N_I(w-1)}$ , then  $\text{Tr}|_{(1, y_t)} < 0$  such that  $(1, y_t)$  is stable but not asymptotically stable, where  $0 < \frac{N_I(1-v_I)(w-1)}{r(w^{N_I}-1)+N_I(w-1)} < 1$ . If  $v_T = v_T^* \wedge 0 \leq v_I < 1 \wedge y_t < \frac{N_I(1-v_I)(w-1)}{r(w^{N_I}-1)+N_I(w-1)}$ , then  $\text{Tr}|_{(1, y_t)} > 0$  such that  $(1, y_t)$  is unstable. If  $v_T = v_T^* \wedge v_I > 1$ , then  $\text{Tr}|_{(1, y_t)} < 0$  such that the entire line of equilibria is stable.

8)  $y_t = 0$ : We find that  $(y_i, 0)$ , where  $0 < y_i < 1$ , is a line of equilibria if and only if  $v_I = 1$ . In this case, the Jacobian at  $(y_i, 0)$  is given by

$$J_{(y_i, 0)} = \begin{pmatrix} 0 & -\frac{(y_i-1)(r((w-1)y_i+1)^{N_I}-1)+N_I(w-1)y_i}{N_I(w-1)} \\ 0 & \frac{(r-1)((w-1)y_i+1)^{N_I}-1}{N_T(w-1)} + v_T \end{pmatrix}. \quad (\text{A.45})$$

We obtain

$$\text{Det}|_{(y_i, 0)} = 0, \quad (\text{A.46})$$

$$\text{Tr}|_{(y_i, 0)} = \frac{(r-1)((w-1)y_i+1)^{N_I}-1}{N_T(w-1)} + v_T \quad (\text{A.47})$$

and

$$\frac{\partial}{\partial y_i} \text{Tr}|_{(y_i, 0)} = -\frac{N_I(1-r)[1+(w-1)y_i]^{N_I-1}}{N_T} < 0. \quad (\text{A.48})$$

If  $v_T = 0$ , then  $\text{Tr}|_{(y_i, 0)} > 0$  such that  $(y_i, 0)$  is unstable. If  $0 < v_T < v_T^* \wedge y_i > \frac{(N_T v_T(1-w)+1)^{1/N_I}-1}{w-1} = \frac{d^{1/N_I}-1}{w-1}$ , then  $\text{Tr}|_{(y_i, 0)} < 0$  such that  $(y_i, 0)$  is stable but not asymptotically stable. If  $0 < v_T < v_T^* \wedge y_i < \frac{d^{1/N_I}-1}{w-1}$ , then  $\text{Tr}|_{(y_i, 0)} > 0$  such that  $(y_i, 0)$  is unstable. Note that we obtain  $0 < \frac{d^{1/N_I}-1}{w-1} < 1$  for  $0 < v_T < v_T^*$ . If  $v_T \geq v_T^*$ , then  $\text{Tr}|_{(y_i, 0)} > 0$  such that  $(y_i, 0)$  is unstable.

9)  $y_t = 1$ : There is no equilibrium on the edge  $(y_i, 1)$ . This is because  $\dot{y}_i = (1-y_i)y_i \left( \frac{r\{1-[1+(w-1)y_i]^{N_I}\}}{N_I(1-w)y_i} + v_I \right) > 0$ , which follows from the combination of  $\frac{1-[1+(w-1)y_i]^{N_I}}{1-w} > 0$  shown in Appendix B5 and  $0 < y_i < 1$ .

### C. Time Average of $(y_i, y_t)$ and the Payoff over a Cycle for $w = 1$

We need to show  $(\bar{y}_i, \bar{y}_t) = (y_{i,Q}, y_{t,Q})$  for  $w = 1$ , where  $\bar{y}_i = \frac{1}{T} \int_0^T y_i dt$ ,  $\bar{y}_t = \frac{1}{T} \int_0^T y_t dt$ ,  $T$  denotes the period of a cycle and  $\mathbf{Q} = (y_{i,Q}, y_{t,Q})$  is given by Eq. (20). By dividing both sides of Eq. (13) by  $y_i(1-y_i) > 0$  and substituting  $w = 1$ , we obtain  $\frac{\dot{y}_i}{y_i(1-y_i)} = \left(\frac{r}{N_I} + 1\right) y_t - 1 + v_I$ . Averaging both sides of the equation over time yields  $0 = \left(\frac{r}{N_I} + 1\right) \bar{y}_t - 1 + v_I$  since  $\frac{1}{T} \int_0^T \frac{\dot{y}_i}{y_i(1-y_i)} dt = 0$ , which follows from  $y_i(0) = y_i(T)$ . On the other hand, Eq. (20) yields  $\left(\frac{r}{N_I} + 1\right) y_{t,Q} - 1 + v_I = 0$ . Therefore, we obtain  $\bar{y}_t = y_{t,Q}$ . Similarly, we can show  $\bar{y}_i = y_{i,Q}$  by starting with dividing both sides of Eq. (14) by  $y_t(1-y_t) > 0$ .

We need to show  $\frac{1}{T} \int_0^T P dt = P(\mathbf{Q})$ , where  $P(y_i, y_t) = \frac{2rN_I y_i y_t + N_I v_I y_i + N_T v_T y_t - a(N_I v_I + N_T v_T)}{N_I + N_T}$ . Because we have shown  $\bar{y}_i = y_{i,Q}$  and  $\bar{y}_t = y_{t,Q}$  above, we only need to show  $\bar{y}_i \bar{y}_t = \bar{y}_i \bar{y}_t$ . To show this, we note that  $\frac{1}{y_i y_t} \frac{d(y_i y_t)}{dt} = y_i y_t \left[ \frac{N_I(1-r)}{N_T} - r - 1 \right] + y_i \left[ \frac{N_I(r-1)}{N_T} - v_I + 1 \right] + y_t(r - v_T + 1) + v_I + v_T - 1$ . Averaging both sides of the equation over time yields  $0 = \bar{y}_i \bar{y}_t \left( \frac{N_I(1-r)}{N_T} - r - 1 \right) + \bar{y}_i \left( \frac{N_I(r-1)}{N_T} - v_I + 1 \right) + \bar{y}_t(r - v_T + 1) + v_I + v_T - 1$  since  $\frac{1}{T} \int_0^T \frac{1}{y_i y_t} \frac{d(y_i y_t)}{dt} dt = 0$ . Therefore, we use  $\bar{y}_i = y_{i,Q} = \frac{N_I v_T}{N_I(1-r)}$  and  $\bar{y}_t = y_{t,Q} = \frac{1-v_I}{1+r}$  to obtain  $\bar{y}_i \bar{y}_t = \frac{\bar{y}_i(N_I(r-1) + N_T(1-v_I)) + \bar{y}_t N_T(r-v_T+1) + N_T(v_I+v_T-1)}{r(N_I+N_T) - N_I + N_T} = \frac{N_T v_T(1-v_I)}{N_I(1-r)(1+r)} = \bar{y}_i \bar{y}_t$ .

### D. Heteroclinic Cycle for $w > 1$

Assume that  $w > 1$ ,  $0 < v_T < v_T^*$  and  $0 \leq v_I < 1$ . We first show that the heteroclinic cycle  $\mathbf{F}_0 \equiv (0, 0) \rightarrow \mathbf{F}_1 \equiv (0, 1) \rightarrow \mathbf{F}_2 \equiv (1, 1) \rightarrow \mathbf{F}_3 \equiv (1, 0) \rightarrow \mathbf{F}_0$  is attracting, i.e., trajectories converge to it. We obtain  $\lambda_1|_{\mathbf{y}} > 0$  and  $\lambda_2|_{\mathbf{y}} < 0$ , where  $\lambda_1|_{\mathbf{y}}$  and  $\lambda_2|_{\mathbf{y}}$  are eigenvalues of the Jacobian at  $\mathbf{y} \in \{\mathbf{F}_0, \mathbf{F}_1, \mathbf{F}_2, \mathbf{F}_3\}$ . Specifically, we obtain  $\lambda_1|_{\mathbf{F}_0} = v_T$ ,  $\lambda_1|_{\mathbf{F}_1} = r + v_I$ ,  $\lambda_1|_{\mathbf{F}_2} = \frac{(1-r)(w^{N_I-1}) - N_T v_T (w-1)}{N_T (w-1)}$ ,  $\lambda_1|_{\mathbf{F}_3} = 1 - v_I$ ,  $\lambda_2|_{\mathbf{F}_0} = -1 + v_I$ ,  $\lambda_2|_{\mathbf{F}_1} = -v_T$ ,  $\lambda_2|_{\mathbf{F}_2} = -\frac{r(w^{N_I-1}) + N_I v_I (w-1)}{N_I (w-1)}$  and  $\lambda_2|_{\mathbf{F}_3} = -\frac{(1-r)(w^{N_I-1}) - N_T v_T (w-1)}{N_T (w-1)}$ . In other words, each  $\mathbf{y}$  is a saddle point. The heteroclinic cycle  $\mathbf{F}_0 \rightarrow \mathbf{F}_1 \rightarrow \mathbf{F}_2 \rightarrow \mathbf{F}_3 \rightarrow \mathbf{F}_0$  is attracting since  $\rho \equiv \left(\frac{-\lambda_2|_{\mathbf{F}_0}}{\lambda_1|_{\mathbf{F}_0}}\right) \left(\frac{-\lambda_2|_{\mathbf{F}_1}}{\lambda_1|_{\mathbf{F}_1}}\right) \left(\frac{-\lambda_2|_{\mathbf{F}_2}}{\lambda_1|_{\mathbf{F}_2}}\right) \left(\frac{-\lambda_2|_{\mathbf{F}_3}}{\lambda_1|_{\mathbf{F}_3}}\right) = \frac{r(w^{N_I-1}) + N_I v_I (w-1)}{N_I (w-1)(r+v_I)} > 1$ , according to the proof of Lemma 1 of Ref. [53].

We show  $\frac{r(w^{N_I-1}) + N_I v_I (w-1)}{N_I (w-1)(r+v_I)} > 1$  as follows. Using  $w > 1$ , we obtain  $1 + N_I(-1+w) - w^{N_I} < 0$  since  $\frac{\partial}{\partial w} [1 + N_I(-1+w) - w^{N_I}] = N_I(1 - w^{N_I-1}) < 0$  and  $[1 + N_I(-1+w) - w^{N_I}]|_{w=1} = 0$ . We then obtain  $1 + N_I(-1+w) - w^{N_I} < 0 \iff r(1 + N_I(-1+w) - w^{N_I}) < 0 \iff N_I r(w-1) < r(w^{N_I} - 1) \iff N_I r(w-1) + N_I v_I (w-1) < r(w^{N_I} - 1) + N_I v_I (w-1) \iff \frac{r(w^{N_I} - 1) + N_I v_I (w-1)}{N_I (w-1)(r+v_I)} > 1$ .

The time average  $\frac{1}{T} \int_0^T (y_i, y_t) dt$  does not converge, where  $(y_i, y_t) = (y_i(t), y_t(t))$  is a trajectory converging to the heteroclinic cycle. According to Theorem 1 of Ref. [53], instead,  $\frac{1}{T} \int_0^T (y_i, y_t) dt$  asymptotically spirals towards the boundary of a polygon (i.e. a quadrangle)  $\mathbf{A}_0 \mathbf{A}_1 \mathbf{A}_2 \mathbf{A}_3$ , where  $\mathbf{A}_i \equiv \frac{\mathbf{F}_{i+1} + \rho_i + 2\mathbf{F}_{i+2} + \rho_i + 2\rho_i + 3\mathbf{F}_{i+3} + \rho_i + 2\rho_i + 3\rho_i + 4\mathbf{F}_{i+4}}{1 + \rho_i + 2\rho_i + 2\rho_i + 3\rho_i + 2\rho_i + 3\rho_i + 4}$ ,  $\rho_i \equiv \frac{-\lambda_2|_{\mathbf{F}_i-1}}{\lambda_1|_{\mathbf{F}_i}}$  and the indices are counted by modulo 4 (e.g.  $\mathbf{F}_4 = \mathbf{F}_0$ ,  $\rho_5 = \rho_1$ ). Because the points  $\mathbf{A}_i, \mathbf{A}_{i+1}$  (with  $i \in \{1, 2, 3, 4\}$ ) and  $\mathbf{F}_{i+1}$  are collinear,  $\frac{1}{T} \int_0^T (y_i, y_t) dt$  asymptotically moves on a line from  $\mathbf{A}_i$  to  $\mathbf{A}_{i+1}$  in the direction of  $\mathbf{F}_{i+1}$  in each cycle [53].

Although the time averages of  $y_i$  and  $y_t$  do not converge, the time average of the payoff  $\bar{P} = \frac{1}{T} \int_0^T P dt$  converges. According to Lemma 1 of Ref. [53], the time for which the trajectory spends near the saddle points  $\mathbf{F}_i$  asymptotically grows  $\rho$  times larger every cycle, whereas the time required to move from a neighbourhood of one saddle point to that of the next one changes little. Thus, we can neglect the latter in comparison with the former. Then, we obtain

$$\begin{aligned} \bar{P} &= \frac{t_0 P(\mathbf{F}_0) + t_1 P(\mathbf{F}_1) + t_2 P(\mathbf{F}_2) + t_3 P(\mathbf{F}_3)}{t_0 + t_1 + t_2 + t_3} \\ &= \frac{P(\mathbf{F}_0) + \frac{t_1}{t_0} P(\mathbf{F}_1) + \frac{t_2}{t_0} P(\mathbf{F}_2) + \frac{t_3}{t_0} P(\mathbf{F}_3)}{1 + \frac{t_1}{t_0} + \frac{t_2}{t_0} + \frac{t_3}{t_0}} \\ &= \frac{P(\mathbf{F}_0) + \frac{t_1}{t_0} P(\mathbf{F}_1) + \frac{t_1}{t_0} \frac{t_2}{t_1} P(\mathbf{F}_2) + \frac{t_1}{t_0} \frac{t_2}{t_1} \frac{t_3}{t_2} P(\mathbf{F}_3)}{1 + \frac{t_1}{t_0} + \frac{t_1}{t_0} \frac{t_2}{t_1} + \frac{t_1}{t_0} \frac{t_2}{t_1} \frac{t_3}{t_2}} \\ &= \frac{P(\mathbf{F}_0) + \rho_1 P(\mathbf{F}_1) + \rho_1 \rho_2 P(\mathbf{F}_2) + \rho_1 \rho_2 \rho_3 P(\mathbf{F}_3)}{1 + \rho_1 + \rho_1 \rho_2 + \rho_1 \rho_2 \rho_3}, \end{aligned} \quad (\text{A.49})$$

where  $t_i$  denotes the time for which the trajectory spends in an arbitrarily small neighbourhood of  $\mathbf{F}_i$  and we have used  $\frac{t_{i+1}}{t_i} = \rho_{i+1}$  from Lemma 1 of Ref. [53]. Note that  $\bar{P}_{\text{hc}}$  is a convex combination of  $P(\mathbf{F}_0), P(\mathbf{F}_1), P(\mathbf{F}_2)$ , and  $P(\mathbf{F}_3)$ . By substituting Eq. (21) with  $(y_i, y_t) = (0, 0), (0, 1), (1, 0)$  and  $(1, 1)$  in Eq. (A.49), we obtain

$$\bar{P}_{\text{hc}} = \frac{1}{(w-1) \left\{ \frac{(r+1)[N_I(1-r) + rN_T v_T]}{N_T v_T [r(w^{N_I-1}) + N_I(w-1)]} - \frac{r}{w^{N_I-1}} \right\} - a(N_I v_I + N_T v_T)} \frac{1}{N_I + N_T}. \quad (\text{A.50})$$

### E. Optimal Incentives

In this section, we calculate the optimal incentive and payoff when  $w \leq 1$  and when  $(w > 1 \wedge N_I \rightarrow \infty) \vee (w \rightarrow \infty)$ .

1)  $w \leq 1$ : To obtain the optimal payoff, we need to know  $\max\{P^*(0, 0), P^*(1, 0), P^*(1, 1), P^*(\mathbf{Q})\}$ . We obtain  $P^*(1, 1) - P^*(\mathbf{Q}) = \frac{r(w^{N_I-1})}{(N_I + N_T)(w-1)} + \frac{(a-1)N_T}{N_I N_T} \epsilon > 0$ . In addition, we have

$$\begin{aligned} \Delta P_1 &\equiv P^*(1, 1) - P^*(1, 0) \\ &= \frac{[(a+1)r - a](w^{N_I} - 1)}{(N_I + N_T)(w-1)} + \frac{aN_I}{N_I + N_T} + (a-1)\epsilon \\ &> 0 \end{aligned} \quad (\text{A.51})$$

for  $0 < w \leq 1$  because  $\Delta P_1$  is a monotonic function of  $w > 0$ , we have  $\Delta P_1|_{w=0} = \frac{a(N_I+r-1)+r}{N_I+N_T} + (a-1)\epsilon > 0$ , and we have  $\Delta P_1|_{w=1} = \frac{(a+1)rN_I}{N_I+N_T} + (a-1)\epsilon > 0$ . Therefore, it holds true that  $\max\{P^*(0,0), P^*(1,0), P^*(1,1), P^*(\mathbf{Q})\} = \max\{P^*(0,0), P^*(1,1)\}$ .

If  $r > r_1^* = \frac{a-1}{a+1}$ , then

$$\begin{aligned} \Delta P_2 &\equiv P^*(1,1) - P^*(0,0) \\ &= \frac{(a+1)r - a + 1}{N_I + N_T} \frac{w^{N_I} - 1}{w - 1} + (1-a)\epsilon > 0. \end{aligned} \quad (\text{A.52})$$

In this case, the optimal payoff is  $P^*(1,1)$ , and the corresponding optimal incentive is  $v_I = 0 \wedge v_T = v_T^* + \epsilon$ . If  $r < r_1^*$ , then  $\Delta P_2 < 0$ . In this case, the optimal payoff is  $P^*(0,0)$ , and the corresponding optimal incentive is  $v_I = 0 \wedge v_T = 0$ .

2)  $(w > 1 \wedge N_I \rightarrow \infty) \vee (w \rightarrow \infty)$ : As  $N_I \rightarrow \infty$ , we obtain

$$\Delta P_1 = P^*(1,1) - P^*(1,0) \rightarrow \frac{(a+1)r - a}{N_I + N_T} \frac{w^{N_I} - 1}{w - 1}. \quad (\text{A.53})$$

The sign of  $\Delta P_1$  is determined by that of  $(a+1)r - a$  since  $\frac{w^{N_I} - 1}{w - 1} > 0$ . Therefore, if  $r > r_2^* = \frac{a}{a+1}$ , then  $P^*(1,1) > P^*(1,0)$ , and if  $r < r_2^*$ , then  $P^*(1,1) < P^*(1,0)$ . As  $N_I \rightarrow \infty$ , we also obtain

$$P^*(1,0) - \bar{P}_{\text{hc}}^* \rightarrow \infty, \quad (\text{A.54})$$

where we remind that  $\bar{P}_{\text{hc}}^*$  denotes the maximum of  $\bar{P}_{\text{hc}}$  with respect to  $v_I$  and  $v_T$ . We prove Eq. (A.54) in Appendix F.

Equations (A.53) and (A.54) imply the following. First, if  $r > r_2^*$ , then  $\max\{P^*(0,0), P^*(1,1), P^*(1,0), \bar{P}_{\text{hc}}^*\} = \max\{P^*(0,0), P^*(1,1)\}$ . Since  $r > r_1^*$ , which follows from  $r_2^* > r_1^*$ , we obtain  $\Delta P_2 = P^*(1,1) - P^*(0,0) > 0$ , which we showed in Eq. (A.52). Therefore,  $\max\{P^*(0,0), P^*(1,1)\} = P^*(1,1)$ ;  $P^*(1,1)$  is the optimal payoff, and the associated optimal incentive is  $v_I = 0 \wedge v_T = v_T^* + \epsilon$ . Second, if  $r < r_2^*$ , then  $\max\{P^*(0,0), P^*(1,1), P^*(1,0), \bar{P}_{\text{hc}}^*\} = \max\{P^*(0,0), P^*(1,0)\}$ . In this case, we obtain  $\Delta P_3 \equiv P^*(1,0) - P^*(0,0) = \frac{1}{N_I + N_T} \left( -N_I a + \frac{1-w^{N_I}}{1-w} \right) - \frac{(a-1)N_I}{N_I + N_T} \epsilon \rightarrow \infty$ . Therefore,  $P^*(1,0)$  is the optimal payoff, and the associated optimal incentive is  $v_I = 1 + \epsilon \wedge v_T = 0$ .

Finally, when  $N_I$  is finite and  $w \rightarrow \infty$ , we have the same outcome via similar calculations.

*F. Proof of Eq. (A.54)*

To prove Eq. (A.54), we first show that  $\bar{P}_{\text{hc}}$  is monotonic or has a local maximum as a function of  $v_T \in (0, v_T^*)$ . Since the denominator of  $\frac{\partial \bar{P}_{\text{hc}}}{\partial v_T}$  is  $(w-1) \{N_I [(1-r^2)(w^{N_I} - 1) - N_T r v_T (w-1)] + N_T r v_T (w^{N_I} - 1)\}^2 > 0$ , the sign of  $\frac{\partial \bar{P}_{\text{hc}}}{\partial v_T}$  is determined by that of its numerator,  $c(v_T) \equiv -v_T^2 N_T^3 a r^2 (w-1) (w^{N_I} - N_I w + N_I - 1)^2 + 2v_T N_I N_T^2 a r (w-1) (w^{N_I} - 1) (r^2 - 1) [(w^{N_I} - 1) - N_I (w-1)] + N_T N_I (1-r^2) (w^{N_I} - 1)^2 [N_I (w-1) (a r^2 - a + 1) + r (w^{N_I} - 1)]$ , which is a quadratic equation of  $v_T$ . Of the two real solutions of  $c(v_T) = 0$ ,

we consider only the larger one,  $v_T = v_T^{\text{hc}} = \frac{\sqrt{a N_I (1-r^2)(w-1)[r(w^{N_I} - 1) + N_I (w-1)] - a N_I (1-r^2)(w-1)}}{a N_T r (w-1)(w^{N_I} - N_I w + N_I - 1)} \times (w^{N_I} - 1)$  because the smaller one is guaranteed to be always negative and  $v_T \geq 0$ . Since the coefficient of the quadratic term  $v_T^2$  is negative, the sign of  $c(v_T)$  over the domain  $(0, v_T^*)$  can be entirely positive, entirely negative, or change from positive to negative just once. In other words,  $\bar{P}_{\text{hc}}$  is monotonic or has a local maximum as a function of  $v_T \in (0, v_T^*)$ . The local maximum of  $\bar{P}_{\text{hc}}$ , if it exists, is realized at  $v_T = v_T^{\text{hc}}$ . Because  $\frac{\partial \bar{P}_{\text{hc}}}{\partial v_I} = -\frac{a N_I}{N_I + N_T} < 0$  implies that the maximum of  $\bar{P}_{\text{hc}}$  in terms of  $v_I$  is realized at  $v_I = 0$  regardless of the value of  $v_T$ , we conclude that  $\bar{P}_{\text{hc}}^*$  is equal to either  $\bar{P}_{\text{hc}}(0 + \epsilon)$ ,  $\bar{P}_{\text{hc}}(v_T^{\text{hc}})$  or  $\bar{P}_{\text{hc}}(v_T^* - \epsilon)$ .

We obtain  $P^*(1,0) - \bar{P}_{\text{hc}}(0 + \epsilon) = \frac{1-w^{N_I}}{(1-w)(N_I+N_T)} - \frac{a N_I}{N_I+N_T} + \frac{[(1-a)N_I + a N_T] \epsilon}{N_I+N_T} \rightarrow \infty$  as  $N_I \rightarrow \infty$  and  $\epsilon \rightarrow 0$ . At  $v_T = v_T^{\text{hc}}$ , we have  $P^*(1,0) - \bar{P}_{\text{hc}}(v_T^{\text{hc}}) = \frac{a N_I^2 r (w-1)^2 + 2(w^{N_I} - 1) \sqrt{a N_I (1-r^2)(w-1)[r(w^{N_I} - 1) + N_I (w-1)]}}{r(w-1)(N_I+N_T)(w^{N_I} - N_I w + N_I - 1)} + \frac{N_I \{a[(r-1)r-1-r-1](w^{N_I} - 1)\}}{r(N_I+N_T)(w^{N_I} - N_I w + N_I - 1)} + \frac{a N_I}{N_I+N_T} \epsilon \rightarrow \infty$  as  $N_I \rightarrow \infty$  and  $\epsilon \rightarrow 0$ . Finally, we have  $P^*(1,0) - \bar{P}_{\text{hc}}(v_T^* - \epsilon) = \frac{[a(1-r)+1](w^{N_I} - 1)}{(w-1)(N_I+N_T)} - \frac{a N_I}{N_I+N_T} + \frac{(1-a)N_I - a N_T}{N_I+N_T} \epsilon \rightarrow \infty$  as  $N_I \rightarrow \infty$  and  $\epsilon \rightarrow 0$ . This concludes the proof of Eq. (A.54).

## REFERENCES

- [1] J. Li, C. Zhang, Q. Sun, Z. Chen, and J. Zhang, "Changing the intensity of interaction based on individual behavior in the iterated prisoner's dilemma game," *IEEE Transactions on Evolutionary Computation*, vol. 21, no. 4, pp. 506–517, 2017.
- [2] R. Chiong and M. Kirley, "Effects of iterated interactions in multi-player spatial evolutionary games," *IEEE Transactions on Evolutionary Computation*, vol. 16, no. 4, pp. 537–555, 2012.
- [3] M. A. Nowak, "Five rules for the evolution of cooperation," *Science*, vol. 314, no. 5805, pp. 1560–1563, 12 2006.
- [4] C. Hauert, F. Michor, M. A. Nowak, and M. Doebeli, "Synergy and discounting of cooperation in social dilemmas," *Journal of Theoretical Biology*, vol. 239, no. 2, pp. 195 – 202, 2006.
- [5] S. Van Segbroeck, F. C. Santos, T. Lenaerts, and J. M. Pacheco, "Reacting differently to adverse ties promotes cooperation in social networks," *Physical Review Letters*, vol. 102, no. 5, pp. 058 105–, 02 2009.
- [6] T. Sasaki, Å. Brännström, U. Dieckmann, and K. Sigmund, "The take-it-or-leave-it option allows small penalties to overcome social dilemmas," *Proceedings of the National Academy of Sciences*, vol. 109, no. 4, pp. 1165–1169, 01 2012.
- [7] J. M. Smith and G. R. Price, "The logic of animal conflict," *Nature*, vol. 246, no. 5427, pp. 15–18, 1973.

- [8] P. D. Taylor and L. B. Jonker, "Evolutionary stable strategies and game dynamics," *Mathematical Biosciences*, vol. 40, no. 1, pp. 145–156, 1978.
- [9] E. D. Bolluyt and C. Comaniciu, "Dynamic influence on replicator evolution for the propagation of competing technologies," *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 5, pp. 899–903, 2019.
- [10] Q. Long, X. Tao, Y. Shi, and S. Zhang, "Evolutionary game analysis among three green-sensitive parties in green supply chains," *IEEE Transactions on Evolutionary Computation*, vol. 25, no. 3, pp. 508–523, 2021.
- [11] G. Bravo and L. Tamburino, "The evolution of trust in non-simultaneous exchange situations," *Rationality and Society*, vol. 20, no. 1, pp. 85–113, 2008.
- [12] N. D. Johnson and A. A. Mislin, "Trust games: A meta-analysis," *Journal of Economic Psychology*, vol. 32, no. 5, pp. 865–889, 2011.
- [13] H. L. J. Ting, X. Kang, T. Li, H. Wang, and C. K. Chu, "On the trust and trust modeling for the future fully-connected digital world: A comprehensive study," *IEEE Access*, vol. 9, pp. 106 743–106 783, 2021.
- [14] D. D. S. Braga, M. Niemann, B. Hellingrath, and F. B. D. L. Neto, "Survey on computational trust and reputation models," *ACM Comput. Surv.*, vol. 51, no. 5, pp. 101:1–101:40, nov 2018.
- [15] J.-H. Cho, K. Chan, and S. Adali, "A survey on trust modeling," *ACM Computing Surveys*, vol. 48, no. 2, pp. 28:1–40, 2015.
- [16] T. Jung, X. Li, W. Huang, Z. Qiao, J. Qian, L. Chen, J. Han, and J. Hou, "Accounttrade: Accountability against dishonest big data buyers and sellers," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 1, pp. 223–234, 2019.
- [17] C. Camerer and K. Weigelt, "Experimental tests of a sequential equilibrium reputation model," *Econometrica*, vol. 56, no. 1, pp. 1–36, 1988.
- [18] J. Berg, J. Dickhaut, and K. McCabe, "Trust, reciprocity, and social history," *Games and Economic Behavior*, vol. 10, no. 1, pp. 122–142, 1995.
- [19] N. Masuda and M. Nakamura, "Coevolution of trustful buyers and cooperative sellers in the trust game," *PLoS one*, vol. 7, no. 9, p. e44169, 2012.
- [20] J. M. McNamara, P. A. Stephens, S. R. X. Dall, and A. I. Houston, "Evolution of trust and trustworthiness: social awareness favours personality differences," *Proceedings of the Royal Society B: Biological Sciences*, vol. 276, no. 1657, pp. 605–613, 02 2009.
- [21] H. Tzieropoulos, "The trust game in neuroscience: A short review," *Social Neuroscience*, vol. 8, no. 5, pp. 407–416, 09 2013.
- [22] I. S. Lim, "Stochastic evolutionary dynamics of trust games with asymmetric parameters," *Physical Review E*, vol. 102, no. 6, pp. 062 419–, 12 2020.
- [23] A. Kumar, V. Capraro, and M. Perc, "The evolution of trust and trustworthiness," *Journal of The Royal Society Interface*, vol. 17, no. 169, p. 20200491, 2020.
- [24] V. Capraro and M. Perc, "Mathematical foundations of moral preferences," *Journal of The Royal Society Interface*, vol. 18, no. 175, p. 20200880, 2021.
- [25] W. Güth and H. Kliemt, "Evolutionarily stable co-operative commitments," *Theory and Decision*, vol. 49, no. 3, pp. 197–222, 2000.
- [26] P. Dasgupta, "Trust as a commodity," in *Trust: Making and Breaking Cooperative Relations*, D. Gambetta, Ed. Department of Sociology, University of Oxford, 2000, ch. 4, pp. 49–72.
- [27] H. Abbass, G. Greenwood, and E. Petraki, "The  $n$ -player trust game and its replicator dynamics," *IEEE Transactions on Evolutionary Computation*, vol. 20, no. 3, pp. 470–474, 2016.
- [28] M. L. Manapat, M. A. Nowak, and D. G. Rand, "Information, irrationality, and the evolution of trust," *Journal of Economic Behavior & Organization*, vol. 90, pp. S57–S75, 2013.
- [29] C. Tarnita, "Fairness and trust in structured populations," *Games*, vol. 6, no. 3, pp. 214–230, 2015.
- [30] I. S. Lim and V. Capraro, "A synergy of institutional incentives and networked structures in evolutionary game dynamics of multiagent systems," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 69, no. 6, pp. 2777–2781, 2022.
- [31] M. Chica, R. Chiong, M. Kirley, and H. Ishibuchi, "A networked  $n$ -player trust game and its evolutionary dynamics," *IEEE Transactions on Evolutionary Computation*, vol. 22, no. 6, pp. 866–878, 2018.
- [32] Z. Hu, X. Li, J. Wang, C. Xia, Z. Wang, and M. Perc, "Adaptive reputation promotes trust in social networks," *IEEE Transactions on Network Science and Engineering*, vol. 8, no. 4, pp. 3087–3098, 2021.
- [33] X. Fang and X. Chen, "Evolutionary dynamics of trust in the  $n$ -player trust game with individual reward and punishment," *The European Physical Journal B*, vol. 94, no. 9, p. 176, 2021.
- [34] S. H. Strogatz, *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering*, 2nd ed. Boca Raton: CRC Press, 2015.
- [35] M. A. Nowak, K. M. Page, and K. Sigmund, "Fairness versus reason in the ultimatum game," *Science*, vol. 289, no. 5485, p. 1773, 09 2000.
- [36] F. Guala and L. Mittone, "Paradigmatic experiments: The dictator game," *The Journal of Socio-Economics*, vol. 39, no. 5, pp. 578–584, 2010.
- [37] J. C. Schank, P. E. Smaldino, and M. L. Miller, "Evolution of fairness in the dictator game by multilevel selection," *Journal of Theoretical Biology*, vol. 382, pp. 64–73, 2015.
- [38] J. E. Snellman, G. Iniguez, J. Kertész, R. A. Barrio, and K. K. Kaski, "Status maximization as a source of fairness in a networked dictator game," *Journal of Complex Networks*, vol. 7, no. 2, pp. 281–305, 2019.
- [39] J. Hofbauer and K. Sigmund, "Evolutionary game dynamics," *Bulletin of the American Mathematical Society*, vol. 40, no. 4, pp. 479–519, 2003.
- [40] K. Sigmund, *The Calculus of Selfishness*. Princeton: Princeton University Press, 2010.
- [41] N. Masuda, "Evolution via imitation among like-minded individuals," *Journal of Theoretical Biology*, vol. 349, pp. 100–108, 2014.
- [42] M. Archetti and I. Scheuring, "Review: Game theory of public goods in one-shot social dilemmas without assortment," *Journal of Theoretical Biology*, vol. 299, pp. 9–20, 2012.
- [43] S. Wang, L. Liu, and X. Chen, "Incentive strategies for the evolution of cooperation: Analysis and optimization," *Europhysics Letters*, vol. 136, no. 6, p. 68002, 2021.
- [44] M. Perc, J. J. Jordan, D. G. Rand, Z. Wang, S. Boccaletti, and A. Szolnoki, "Statistical physics of human cooperation," *Physics Reports*, vol. 687, pp. 1–51, 2017.
- [45] J. Zhang and M. Cao, "Strategy competition dynamics of multi-agent systems in the framework of evolutionary game theory," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 67, no. 1, pp. 152–156, 2020.
- [46] Y. Dong, T. Sasaki, and B. Zhang, "The competitive advantage of institutional reward," *Proceedings of the Royal Society B: Biological Sciences*, vol. 286, no. 1899, p. 20190001, 2019.
- [47] R. Cressman, J.-W. Song, B.-Y. Zhang, and Y. Tao, "Cooperation and evolutionary dynamics in the public goods game with institutional incentives," *Journal of Theoretical Biology*, vol. 299, pp. 144–151, 2012.
- [48] M. Perc, "Sustainable institutionalized punishment requires elimination of second-order free-riders," *Scientific Reports*, vol. 2, no. 1, p. 344, 2012.
- [49] F. C. Santos, M. D. Santos, and J. M. Pacheco, "Social diversity promotes the emergence of cooperation in public goods games," *Nature*, vol. 454, pp. 213 EP–, 07 2008.
- [50] U. Alvarez-Rodriguez, F. Battiston, G. F. de Arruda, Y. Moreno, M. Perc, and V. Latora, "Evolutionary dynamics of higher-order interactions in social networks," *Nature Human Behaviour*, vol. 5, pp. 586–595, 2021.
- [51] N. Samian, Z. A. Zukarnain, W. K. G. Seah, A. Abdullah, and Z. M. Hanapi, "Cooperation stimulation mechanisms for wireless multihop networks: A survey," *Journal of Network and Computer Applications*, vol. 54, pp. 88–106, 2015.
- [52] B. M. C. Silva, J. J. P. C. Rodrigues, N. Kumar, and G. Han, "Cooperative strategies for challenged networks and applications: A survey," *IEEE Systems Journal*, vol. 11, no. 4, pp. 2749–2760, 2017.
- [53] A. Gaunersdorfer, "Time averages for heteroclinic attractors," *SIAM Journal on Applied Mathematics*, vol. 52, no. 5, pp. 1476–1489, 1992.



**Ik Soo Lim** Ik Soo Lim received a PhD from the Swiss Federal Institute of Technology at Lausanne (EPFL). He is a senior lecturer of computer science at the School of Computing and Mathematical Sciences, University of Greenwich, London. His research interests include multi-agent systems, evolutionary game theory and data visualisation among others. He has publications in journals of physics, biology, psychology and computing.



**Naoki Masuda** Naoki Masuda received his PhD in 2002 from the University of Tokyo. He worked as Lecturer and then Associate Professor at the University of Tokyo, Japan, between 2006 and 2014. Then, he worked as Senior Lecturer and Associate Professor at the University of Bristol, UK, between 2014 and 2019. He moved to Department of Mathematics at University at Buffalo in 2019 as Associate Professor and has been full Professor since 2021. His research interests include network science and mathematical biology.