

PNAS



1

2 **Supplementary Information for**

3 **Evolutionary stability of cooperation in indirect reciprocity with noisy and private reputation**

4 **Yuma Fujimoto, and Hisashi Ohtsuki**

5 **Yuma Fujimoto.**

6 **E-mail: fujimoto_yuma@soken.ac.jp**

7 **This PDF file includes:**

8 Supplementary text

9 Figs. S1 to S4

10 Tables S1 to S2

11 SI References

12 Supporting Information Text

13 1. Calculation of joint distribution of goodnesses

14 This section proposes an analytical method to obtain the reputation structure under indirect reciprocity. We assume a situation
 15 where rare mutants with norm M of ratio δ invade other wild-types with norm W of ratio $1 - \delta$. We denote the population ratio
 16 of norm $A \in \{W, M\}$ as ρ_A ; $\rho_A = 1 - \delta$ when $A = W$, while $\rho_A = \delta$ when $A = M$. To characterize the reputation structure, we
 17 define p_{iA} as a proportion of individuals of norm A who assign good reputations to individual i . We call p_{iA} a goodness of
 18 individual i from norm $A \in \{W, M\}$.

19 In the following, let us consider a stochastic transition of p_{iA} in each round. In a single round, a recipient and a donor are
 20 chosen and labeled as i_R and i_D , respectively. In this round, $p_{i_D A_O}$, i.e., the goodness of donor from norm A_O , changes into
 21 the next goodness $p'_{i_D A_O}$ for all $A_O \in \{W, M\}$. Below, we formulate the stochastic change separately for cases that the donor
 22 chooses to cooperate or defect.

C-map case: First, we consider a case that the donor cooperates with the recipient, occurring with a probability of

$$h(p_{i_R A_D}) := p_{i_R A_D}(1 - e_1) + (1 - p_{i_R A_D})e_1. \quad [1]$$

Here, e_1 is the probability of action error, with which the donor chooses an opposite action to the intended one. In this case,
 $N\rho_{A_O}p'_{i_D A_O}$, i.e., the number of observers with norm A_O who give good reputations to the donor in the next round, follows a
 probability distribution of

$$N\rho_{A_O}p'_{i_D A_O} \sim N_1 + N_2, \quad [2]$$

$$N_1 \sim \mathcal{B}(N\rho_{A_O}p_{i_R A_O}, a_{A_O}^{GC}), \quad [3]$$

$$N_2 \sim \mathcal{B}(N\rho_{A_O}(1 - p_{i_R A_O}), a_{A_O}^{BC}). \quad [4]$$

Here, $\mathcal{B}(n, a)$ denotes a binomial distribution with success probability a and trial number n . In addition, $a_{A_O}^{XY}$ denotes the
 probability that an observer with norm A_O who evaluates the recipient as $X \in \{G, B\}$ newly gives a good reputation to the
 donor whose action is $Y \in \{C, D\}$. $a_{A_O}^{GC}$, $a_{A_O}^{BC}$, $a_{A_O}^{GD}$, and $a_{A_O}^{BD}$ are obtained by converting corresponding G and B pivots into
 1 - e_2 and e_2 in Table 1 of the main manuscript. Instead of Eq. (3), we use a shorthand notation;

$$N\rho_{A_O}p'_{i_D A_O} \sim \mathcal{B}(N\rho_{A_O}p_{i_R A_O}, a_{A_O}^{GC}) + \mathcal{B}(N\rho_{A_O}(1 - p_{i_R A_O}), a_{A_O}^{BC}). \quad [5]$$

Because $N\rho_{A_O}$ is sufficiently large, the mean and variance of $p'_{i_D A_O}$ are given by

$$\mathbb{E}[p'_{i_D A_O}] = p_{i_R A_O} \underbrace{(a_{A_O}^{GC} - a_{A_O}^{BC})}_{=: \Delta f_{A_O}^C} + a_{A_O}^{BC} \quad (=: f_{A_O}^C(p_{i_R A_O})), \quad [6]$$

$$\text{Var}[p'_{i_D A_O}] = \frac{p_{i_R A_O} a_{A_O}^{GC}(1 - a_{A_O}^{GC}) + (1 - p_{i_R A_O}) a_{A_O}^{BC}(1 - a_{A_O}^{BC})}{N\rho_{A_O}} = \frac{e_2(1 - e_2)}{N\rho_{A_O}} \quad (=: \rho_{A_O}^{-1} s^2). \quad [7]$$

23 Here, e_2 is the probability of assessment error, with which an observer assigns an opposite reputation to the intended one. In
 24 Eq. (6), $f_{A_O}^C$ represents a map from the recipient's goodness in the present round to the donor's goodness in the next round.
 25 Because this map is applied only when the donor cooperates, we call it "C-map".

D-map case: On the other hand, we consider a case that the donor defects with the recipient, occurring with a probability
 of

$$1 - h(p_{i_R A_D}) = (1 - p_{i_R A_D})(1 - e_1) + p_{i_R A_D}e_1. \quad [8]$$

In this case, $N\rho_{A_O}p'_{i_D A_O}$ follows a probability distribution of

$$N\rho_{A_O}p'_{i_D A_O} \sim \mathcal{B}(N\rho_{A_O}p_{i_R A_O}, a_{A_O}^{GD}) + \mathcal{B}(N\rho_{A_O}(1 - p_{i_R A_O}), a_{A_O}^{BD}). \quad [9]$$

From this equation, the mean and variance of $p'_{i_D A_O}$ are given by

$$\mathbb{E}[p'_{i_D A_O}] = p_{i_R A_O} \underbrace{(a_{A_O}^{GD} - a_{A_O}^{BD})}_{=: \Delta f_{A_O}^D} + a_{A_O}^{BD} \quad (=: f_{A_O}^D(p_{i_R A_O})), \quad [10]$$

$$\text{Var}[p'_{i_D A_O}] = \frac{p_{i_R A_O} a_{A_O}^{GD}(1 - a_{A_O}^{GD}) + (1 - p_{i_R A_O}) a_{A_O}^{BD}(1 - a_{A_O}^{BD})}{N\rho_{A_O}} = \frac{e_2(1 - e_2)}{N\rho_{A_O}} \quad (=: \rho_{A_O}^{-1} s^2). \quad [11]$$

26 Because the map $f_{A_O}^D$ is applied when the donor defects, we call it D-map in the same way as C-map.

27 The above C-map $f_{S_k}^C$ and D-map $f_{S_k}^D$ are illustrated in Fig. S1 for all $S_k \in \mathcal{S}$.

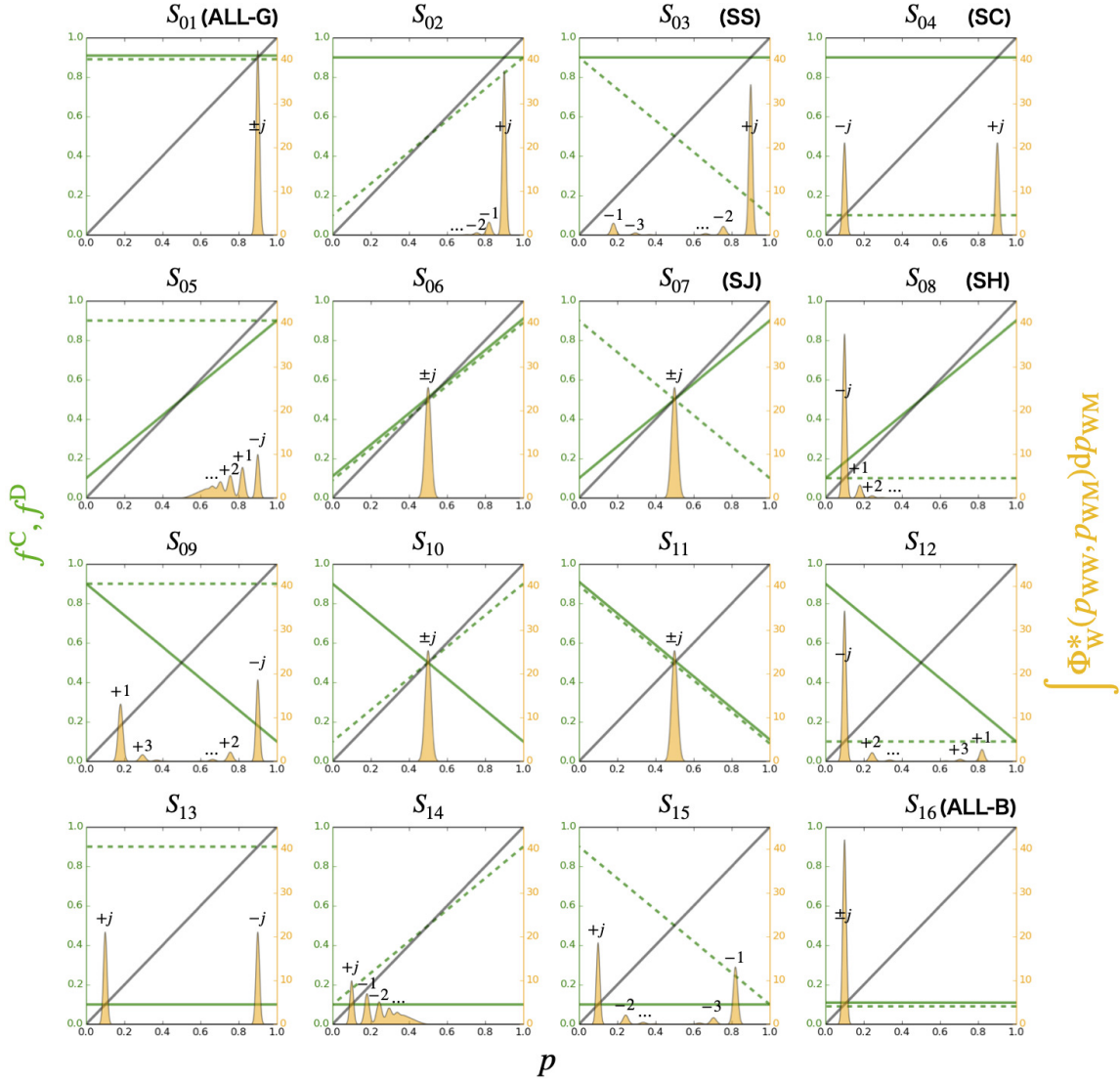


Fig. S1. Materials for the reputation structure for all the second-order norms $W = S_k$. Green solid (resp. dashed) lines indicate C-map f_W^C (resp. D-map f_W^D) of the norm. Gray lines indicate the identity map, which shows the fixed points of the C-map and D-map as the crossing points with these maps. The orange distribution shows the probability density function of goodnesses p_{WW} when $W = S_k$. All the panels are output under $N = 2000$, $\delta = 0$, and $(\epsilon_1, \epsilon_2) = (0, 0.1)$. Numbers over each peak indicate j .

28 2. Time evolution of reputation structure

Because the population of wild-types and mutants are sufficiently large, we can continualize the distribution of individual goodness p_{iA} with separating into the cases that the norm of individual i is W or M. In the following, $p_{AA'}$ denotes a continualized goodness of an individual with norm A in the eyes of individuals with norm A' . Let us consider a time change of the distribution of $p_{AA'}$. As shown above, however, we should keep in mind that p_{AW} and p_{AM} are simultaneously changed by the C-map or D-map. Thus, we consider dynamics of $\Phi_A(p_{AW}, p_{AM})$, a joint probability distribution of p_{AW} and p_{AM} . Note that a norm of the chosen recipient is M only with a probability of δ , which contributes the dynamics of Φ_A only in a scale of $O(\delta)$. By ignoring this scale of $O(\delta)$, the dynamics of Φ_A is given by

$$\begin{aligned} \frac{d}{dt} \Phi_A(p_{AW}, p_{AM}) = & -\Phi_A(p_{AW}, p_{AM}) + \int_0^1 \int_0^1 \{h(p'_{WA})g(p_{AW}; f_W^C(p'_{WW}), \rho_W^{-1} s^2)g(p_{AM}; f_M^C(p'_{WM}), \rho_M^{-1} s^2) \\ & + (1 - h(p'_{WA}))g(p_{AW}; f_W^D(p'_{WW}), \rho_W^{-1} s^2)g(p_{AM}; f_M^D(p'_{WM}), \rho_M^{-1} s^2)\} \\ & \times \Phi_W(p'_{WW}, p'_{WM}) dp'_{WW} dp'_{WM}. \end{aligned} \quad [12]$$

Here, $g(p; \mu, \sigma^2)$ denotes a Gaussian function with the mean μ and variance σ^2 as

$$g(p; \mu, \sigma^2) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(p - \mu)^2}{2\sigma^2}\right). \quad [13]$$

Equation (12) explains an update of the donor's goodness per time. The first (resp. second) term on the right side represents decrements (increments) by updating goodnnesses. In detail, $\Phi_W(p'_{WW}, p'_{WM})$ in the second term shows the density that the recipient's goodness is p'_{WW} (resp. p'_{WM}) in the eyes of wild-types (resp. mutants). $h(p'_{WA})$ shows the probability that the donor cooperates, and the donor's goodnnesses after the update in the eyes of observers with norm W and M are described by $g(p_{AW}; f_W^C(p'_{WW}), \rho_W^{-1} s^2)$ and $g(p_{AM}; f_M^C(p'_{WM}), \rho_M^{-1} s^2)$, respectively. A similar explanation holds when the donor chooses to defect.

The equilibrium state of Eq. (12), i.e., Φ_A^* , satisfies

$$\begin{aligned} \Phi_A^*(p_{AW}, p_{AM}) = & \int_0^1 \int_0^1 \{h(p'_{WA})g(p_{AW}; f_W^C(p'_{WW}), \rho_W^{-1} s^2)g(p_{AM}; f_M^C(p'_{WM}), \rho_M^{-1} s^2) \\ & + (1 - h(p'_{WA}))g(p_{AW}; f_W^D(p'_{WW}), \rho_W^{-1} s^2)g(p_{AM}; f_M^D(p'_{WM}), \rho_M^{-1} s^2)\} \\ & \times \Phi_W^*(p'_{WW}, p'_{WM}) dp'_{WW} dp'_{WM}. \end{aligned} \quad [14]$$

To solve this equation, we assume that the equilibrium state can be described by a summation of two-dimensional Gaussian functions without correlation as

$$\Phi_A^*(p_{AW}, p_{AM}) = \sum_j q_{Aj} g(p_{AW}; \mu_{AWj}, \rho_W^{-1} \sigma_{AWj}^2) g(p_{AM}; \mu_{AMj}, \rho_M^{-1} \sigma_{AMj}^2). \quad [15]$$

The assumption of Gaussian is justified by the above transition process of the donor's goodness, where the goodness is virtually determined only by the mean and variance in a sufficiently large population. No correlation is assumed because the variance is given independently by observers with different norms.

We now derive equations in which the equilibrium state satisfies for each norm $A \in \{W, M\}$. First, substituting Eq. (15)

into Eq. (14) for $A = W$, we obtain

$$\begin{aligned}
& \sum_j q_{Wj} g(p_{WW}; \mu_{WWj}, \rho_W^{-1} \sigma_{WWj}^2) g(p_{WM}; \mu_{WMj}, \rho_M^{-1} \sigma_{WMj}^2) \\
&= \int_0^1 \int_0^1 \{h(p'_{WW}) g(p_{WW}; f_W^C(p'_{WW}), \rho_W^{-1} s^2) g(p_{WM}; f_M^C(p'_{WM}), \rho_M^{-1} s^2) \\
&\quad + (1 - h(p'_{WW})) g(p_{WW}; f_W^D(p'_{WW}), \rho_W^{-1} s^2) g(p_{WM}; f_M^D(p'_{WM}), \rho_M^{-1} s^2)\} \\
&\quad \times \sum_j q_{Wj} g(p'_{WW}; \mu_{WWj}, \rho_W^{-1} \sigma_{WWj}^2) g(p'_{WM}; \mu_{WMj}, \rho_M^{-1} \sigma_{WMj}^2) dp'_{WW} dp'_{WM}, \\
&= \sum_j q_{Wj} \int_0^1 \int_0^1 \{h(p'_{WW}) g(p_{WW}; f_W^C(p'_{WW}), \rho_W^{-1} s^2) g(p_{WM}; f_M^C(p'_{WM}), \rho_M^{-1} s^2) \\
&\quad + (1 - h(p'_{WW})) g(p_{WW}; f_W^D(p'_{WW}), \rho_W^{-1} s^2) g(p_{WM}; f_M^D(p'_{WM}), \rho_M^{-1} s^2)\} \\
&\quad \times g(p'_{WW}; \mu_{WWj}, \rho_W^{-1} \sigma_{WWj}^2) g(p'_{WM}; \mu_{WMj}, \rho_M^{-1} \sigma_{WMj}^2) dp'_{WW} dp'_{WM}, \\
&\simeq \sum_j q_{Wj} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{h(\mu_{WWj}) g(p_{WW}; f_W^C(\mu_{WWj}), \rho_W^{-1} s^2) g(p_{WM}; f_M^C(\mu_{WMj}), \rho_M^{-1} s^2) \\
&\quad + (1 - h(\mu_{WWj})) g(p_{WW}; f_W^D(\mu_{WWj}), \rho_W^{-1} s^2) g(p_{WM}; f_M^D(\mu_{WMj}), \rho_M^{-1} s^2)\} \\
&\quad \times g(\mu_{WWj}; \mu_{WWj}, \rho_W^{-1} \sigma_{WWj}^2) g(\mu_{WMj}; \mu_{WMj}, \rho_M^{-1} \sigma_{WMj}^2) d\mu_{WWj} d\mu_{WMj}, \\
&= \sum_j q_{Wj} \{h(\mu_{WWj}) g(p_{WW}; f_W^C(\mu_{WWj}), \rho_W^{-1} (s^2 + (\Delta f_W^C)^2 \sigma_{WWj}^2)) g(p_{WM}; f_M^C(\mu_{WMj}), \rho_M^{-1} (s^2 + (\Delta f_M^C)^2 \sigma_{WMj}^2))\} \\
&\quad + (1 - h(\mu_{WWj})) g(p_{WW}; f_W^D(\mu_{WWj}), \rho_W^{-1} (s^2 + (\Delta f_W^D)^2 \sigma_{WWj}^2)) g(p_{WM}; f_M^D(\mu_{WMj}), \rho_M^{-1} (s^2 + (\Delta f_M^D)^2 \sigma_{WMj}^2))\}. \quad [16]
\end{aligned}$$

This equation gives a constraint for $(q_{Wj}, \mu_{WWj}, \sigma_{WWj}^2, \mu_{WMj}, \sigma_{WMj}^2)$. Next, when $A = M$, in a similar manner, we obtain

$$\begin{aligned}
& \sum_j q_{Mj} g(p_{MW}; \mu_{MWj}, \rho_W^{-1} \sigma_{MWj}^2) g(p_{MM}; \mu_{MMj}, \rho_M^{-1} \sigma_{MMj}^2) \\
&= \int_0^1 \int_0^1 \{h(p'_{MW}) g(p_{MW}; f_W^C(p'_{MW}), \rho_W^{-1} s^2) g(p_{MM}; f_M^C(p'_{MM}), \rho_M^{-1} s^2) \\
&\quad + (1 - h(p'_{MW})) g(p_{MW}; f_W^D(p'_{MW}), \rho_W^{-1} s^2) g(p_{MM}; f_M^D(p'_{MM}), \rho_M^{-1} s^2)\} \\
&\quad \times \sum_j q_{Mj} g(p'_{MW}; \mu_{MWj}, \rho_W^{-1} \sigma_{MWj}^2) g(p'_{MM}; \mu_{MMj}, \rho_M^{-1} \sigma_{MMj}^2) dp'_{MW} dp'_{MM} \\
&= \sum_j q_{Mj} \int_0^1 \int_0^1 \{h(p'_{MW}) g(p_{MW}; f_W^C(p'_{MW}), \rho_W^{-1} s^2) g(p_{MM}; f_M^C(p'_{MM}), \rho_M^{-1} s^2) \\
&\quad + (1 - h(p'_{MW})) g(p_{MW}; f_W^D(p'_{MW}), \rho_W^{-1} s^2) g(p_{MM}; f_M^D(p'_{MM}), \rho_M^{-1} s^2)\} \\
&\quad \times \sum_j q_{Mj} g(p'_{MW}; \mu_{MWj}, \rho_W^{-1} \sigma_{MWj}^2) g(p'_{MM}; \mu_{MMj}, \rho_M^{-1} \sigma_{MMj}^2) dp'_{MW} dp'_{MM} \\
&\simeq \sum_j q_{Mj} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{h(\mu_{MWj}) g(p_{MW}; f_W^C(\mu_{MWj}), \rho_W^{-1} s^2) g(p_{MM}; f_M^C(\mu_{MMj}), \rho_M^{-1} s^2) \\
&\quad + (1 - h(\mu_{MWj})) g(p_{MW}; f_W^D(\mu_{MWj}), \rho_W^{-1} s^2) g(p_{MM}; f_M^D(\mu_{MMj}), \rho_M^{-1} s^2)\} \\
&\quad \times \sum_j q_{Mj} g(\mu_{MWj}; \mu_{MWj}, \rho_W^{-1} \sigma_{MWj}^2) g(\mu_{MMj}; \mu_{MMj}, \rho_M^{-1} \sigma_{MMj}^2) d\mu_{MWj} d\mu_{MMj} \\
&= \sum_j q_{Mj} \{h(\mu_{MWj}) g(p_{MW}; f_W^C(\mu_{MWj}), \rho_W^{-1} (s^2 + (\Delta f_W^C)^2 \sigma_{MWj}^2)) g(p_{MM}; f_M^C(\mu_{MMj}), \rho_M^{-1} (s^2 + (\Delta f_M^C)^2 \sigma_{MMj}^2))\} \\
&\quad + (1 - h(\mu_{MWj})) g(p_{MW}; f_W^D(\mu_{MWj}), \rho_W^{-1} (s^2 + (\Delta f_W^D)^2 \sigma_{MWj}^2)) g(p_{MM}; f_M^D(\mu_{MMj}), \rho_M^{-1} (s^2 + (\Delta f_M^D)^2 \sigma_{MMj}^2))\}. \quad [17]
\end{aligned}$$

38 This equation gives a constraint for $(q_{Mj}, \mu_{MWj}, \sigma_{MWj}^2, \mu_{MMj}, \sigma_{MMj}^2)$.

To solve Eq. (16), let us consider a set of solutions $\{(\mu_{WWj}, \mu_{WMj})\}_j$. From the equilibrium condition of Eq. (16), the equal set must be restored by applying C-map and D-map to all the elements of the set. In other words, the condition is given by

$$\{(\mu_{WWj}, \mu_{WMj})\}_j = \{(f_W^C(\mu_{WWj}), f_M^C(\mu_{WMj}))\}_j \cup \{(f_W^D(\mu_{WWj}), f_M^D(\mu_{WMj}))\}_j. \quad [18]$$

Similarly, to obtain the equilibrium condition for Eq. (17), we should consider a set of solutions $\{(\mu_{MWj}, \mu_{MMj})\}_j$ satisfying

$$\{(\mu_{MWj}, \mu_{MMj})\}_j = \{(f_W^C(\mu_{MWj}), f_M^C(\mu_{MMj}))\}_j \cup \{(f_W^D(\mu_{MWj}), f_M^D(\mu_{MMj}))\}_j. \quad [19]$$

Here, although the appearance of the variables is different, the problems are essentially between Eq. (18) and Eq. (19). Thus, the problem to be solved is

$$\{(\mu_{j,W}, \mu_{j,M})\}_j = \{(f_W^C(\mu_{j,W}), f_M^C(\mu_{j,M}))\}_j \cup \{(f_W^D(\mu_{j,W}), f_M^D(\mu_{j,M}))\}_j. \quad [20]$$

Furthermore, because $W, M \in \mathcal{S}$, we can generalize the problem as

$$\{(\mu_{j,S_{01}}, \dots, \mu_{j,S_{16}})\}_j = \{(f_{S_{01}}^C(\mu_{j,S_{01}}), \dots, f_{S_{16}}^C(\mu_{j,S_{16}}))\}_j \cup \{(f_{S_{01}}^D(\mu_{j,S_{01}}), \dots, f_{S_{16}}^D(\mu_{j,S_{16}}))\}_j. \quad [21]$$

Now, for all $S_k \in \mathcal{S}$ and, let us consider a set $\{\mu_{j,S_k}\}_{j \in \mathbb{Z} \setminus \{0\}}$ satisfying

$$\mu_{+1,S_k} = f_{S_k}^C(\mu_{-1,S_k}) = f_{S_k}^C(\mu_{-2,S_k}) = \dots, \quad [22]$$

$$\mu_{+(j+1),S_k} = f_{S_k}^C(\mu_{+j,S_k}), \quad [23]$$

$$\mu_{-1,S_k} = f_{S_k}^D(\mu_{+1,S_k}) = f_{S_k}^D(\mu_{+2,S_k}) = \dots, \quad [24]$$

$$\mu_{-(j+1),S_k} = f_{S_k}^D(\mu_{-j,S_k}), \quad [25]$$

(the proof for these equations will be given later). This set $\{\mu_{j,S_k}\}_{j \in \mathbb{Z} \setminus \{0\}}$ gives a solution to problem Eq. (21), and thus solves Eq. (16) and Eq. (17). In Eq. (22)-Eq. (25), we consistently label each μ_{j,S_k} such that sequentially applying C-map (resp. D-map) $j(> 0)$ times leads to label $+j$ (resp. $-j$) (see the illustration in Fig. S1). In the following, we show that such $\{\mu_{j,S_k}\}_{j \in \mathbb{Z} \setminus \{0\}}$ actually exists for all k .

1. When neither C-map nor D-map is constant: First, we consider a case of $\Delta f_{S_k}^C \neq 0$ and $\Delta f_{S_k}^D \neq 0$. Four norms of $k = 06, 07, 09, 10$ correspond to this case. In this case, C-map and D-map have the same fixed point (at $1/2$). Only the position given by this fixed point is achieved at equilibrium. Indeed, if we substitute $\mu_{j,S_k} = 1/2$ for all $j \in \mathbb{Z} \setminus \{0\}$, Eq. (22)-Eq. (25) are simultaneously satisfied without any contradiction.

2. When both C-map and D-map are constant: Second, we consider a case of $\Delta f_{S_k}^C = 0$ and $\Delta f_{S_k}^D = 0$. Four norms of $k = 01, 04, 13, 16$ correspond to this case. Because $f_{S_k}^C$ is a constant map, Eq. (22) and Eq. (23) are satisfied by substituting the mapped value of this map into μ_{+j,S_k} for all $j = 1, 2, \dots$. In the same way, because $f_{S_k}^D$ is a constant map, Eq. (24) and Eq. (25) are satisfied by substituting the mapped value into μ_{-j,S_k} for all $j = 1, 2, \dots$. Thus, no contradiction occurs.

3. When only C-map is constant: Third, we consider a case of $\Delta f_{S_k}^C = 0$ and $\Delta f_{S_k}^D \neq 0$. Four norms $k = 02, 03, 14, 15$ correspond to this case. Because $f_{S_k}^C$ is a constant map, Eq. (22) and Eq. (23) are satisfied by substituting the mapped value of this map into μ_{+j,S_k} for all $j = 1, 2, \dots$. Then, we define μ_{-1,S_k} as the value to which D-map maps all the same value $\mu_{+1,S_k} = \mu_{+2,S_k} = \dots$, and Eq. (24) is satisfied. Finally, we sequentially define $\mu_{-2,S_k}, \mu_{-3,S_k}, \dots$ by applying D-map to μ_{-1,S_k} one by one. Thus, no contradiction occurs.

4. When only D-map is constant: Finally, we consider a case of $\Delta f_{S_k}^C \neq 0$ and $\Delta f_{S_k}^D = 0$. Four norms $k = 05, 08, 09, 12$ correspond to this case. Because $f_{S_k}^D$ is a constant map, Eq. (24), Eq. (25) are satisfied by substituting the mapped value of this map into μ_{-j,S_k} for all $j = 1, 2, \dots$. Then, we define μ_{+1,S_k} as the value to which D-map maps all the same value $\mu_{-1,S_k} = \mu_{-2,S_k} = \dots$, and Eq. (24) is satisfied. Finally, we sequentially define $\mu_{+2,S_k}, \mu_{+3,S_k}, \dots$ by applying C-map to μ_{+1,S_k} one by one. Thus, no contradiction occurs.

As summarized in Table S1, the set $\{\mu_{j,S_k}\}_{j \in \mathbb{Z} \setminus \{0\}}$ can be analytically described. Furthermore, we also define σ_{j,S_k}^2 as the variance in Gaussian corresponding to the mean μ_{j,S_k} . Similarly to the mean values above, we solve the variances as

$$(\sigma_{WWj}^2, \sigma_{WMj}^2) = (\sigma_{MWj}^2, \sigma_{MMj}^2) = (\sigma_{j,W}^2, \sigma_{j,M}^2). \quad [26]$$

The recursion that the set $\{\sigma_{j,S_k}^2\}_{j \in \mathbb{Z} \setminus \{0\}}$ should satisfy is

$$\sigma_{+1,S_k}^2 = s^2 + (\Delta f_{S_k}^C)^2 \sigma_{-1,S_k}^2 = s^2 + (\Delta f_{S_k}^C)^2 \sigma_{-2,S_k}^2 = \dots, \quad [27]$$

$$\sigma_{+(j+1),S_k}^2 = s^2 + (\Delta f_{S_k}^C)^2 \sigma_{+j,S_k}^2, \quad [28]$$

$$\sigma_{-1,S_k}^2 = s^2 + (\Delta f_{S_k}^D)^2 \sigma_{+1,S_k}^2 = s^2 + (\Delta f_{S_k}^D)^2 \sigma_{+2,S_k}^2 = \dots, \quad [29]$$

$$\sigma_{-(j+1),S_k}^2 = s^2 + (\Delta f_{S_k}^D)^2 \sigma_{-j,S_k}^2, \quad [30]$$

and the solution exists for all S_k (see Table. S1 for the solution of these equations). Table. S1 shows $\{(\mu_{j,S_k}, \sigma_{j,S_k}^2)\}_{j \in \mathbb{Z} \setminus \{0\}}$.

S_k	μ_{+j, S_k}	μ_{-j, S_k}	σ_{+j, S_k}^2	σ_{-j, S_k}^2
S_{01}	$1 - e_2$	$1 - e_2$	$\frac{e_2(1 - e_2)}{N}$	$\frac{e_2(1 - e_2)}{N}$
S_{02}	$1 - e_2$	$\frac{1 + (1 - 2e_2)^{j+1}}{2}$	$\frac{e_2(1 - e_2)}{N}$	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$
S_{03}	$1 - e_2$	$\frac{1 - \{(1 - 2e_2)\}^{j+1}}{2}$	$\frac{e_2(1 - e_2)}{N}$	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$
S_{04}	$1 - e_2$	e_2	$\frac{e_2(1 - e_2)}{N}$	$\frac{e_2(1 - e_2)}{N}$
S_{05}	$\frac{1 + (1 - 2e_2)^{j+1}}{2}$	$1 - e_2$	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$	$\frac{e_2(1 - e_2)}{N}$
S_{07}	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{4N}$	$\frac{1}{4N}$
S_{08}	$\frac{1 - (1 - 2e_2)^{j+1}}{2}$	e_2	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$	$\frac{e_2(1 - e_2)}{N}$
S_{09}	$\frac{1 - \{(1 - 2e_2)\}^{j+1}}{2}$	$1 - e_2$	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$	$\frac{e_2(1 - e_2)}{N}$
S_{12}	$\frac{1 + \{(1 - 2e_2)\}^{j+1}}{2}$	e_2	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$	$\frac{e_2(1 - e_2)}{N}$
S_{13}	e_2	$1 - e_2$	$\frac{e_2(1 - e_2)}{N}$	$\frac{e_2(1 - e_2)}{N}$
S_{14}	e_2	$\frac{1 - (1 - 2e_2)^{j+1}}{2}$	$\frac{e_2(1 - e_2)}{N}$	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$
S_{15}	e_2	$\frac{1 + \{(1 - 2e_2)\}^{j+1}}{2}$	$\frac{e_2(1 - e_2)}{N}$	$\frac{1 - (1 - 2e_2)^{2(j+1)}}{4N}$
S_{16}	e_2	e_2	$\frac{e_2(1 - e_2)}{N}$	$\frac{e_2(1 - e_2)}{N}$

Table S1. Analytical solution of Gaussian functions. We omit S_{06} , S_{10} , and S_{11} because the results are identical to those of S_{07} (SJ).

We also calculate a set of the masses of Gaussian functions, i.e., $\{q_{Wj}\}_{j \in \mathbb{Z} \setminus \{0\}}$ and $\{q_{Mj}\}_{j \in \mathbb{Z} \setminus \{0\}}$. By substituting the values in Table S1 into Eq. (16), we obtain the following relational expressions

$$q_{W+1} = \sum_{j=1}^{\infty} h(\mu_{-j, W}) q_{W-j}, \quad [31]$$

$$q_{W+j} = h(\mu_{+(j-1), W}) q_{W+(j-1)} \quad (j = 2, \dots, \infty), \quad [32]$$

$$q_{W-1} = \sum_{j=1}^{\infty} (1 - h(\mu_{+j, W})) q_{W+j}, \quad [33]$$

$$q_{W-j} = (1 - h(\mu_{-(j-1), W})) q_{W-(j-1)} \quad (j = 2, \dots, \infty). \quad [34]$$

Similarly, by substituting the values in Table S1 into Eq. (17), we obtain

$$q_{M+1} = \sum_{j=1}^{\infty} h(\mu_{-j, M}) q_{W-j}, \quad [35]$$

$$q_{M+j} = h(\mu_{+(j-1), M}) q_{W+(j-1)} \quad (j = 2, \dots, \infty), \quad [36]$$

$$q_{M-1} = \sum_{j=1}^{\infty} (1 - h(\mu_{+j, M})) q_{W+j}, \quad [37]$$

$$q_{M-j} = (1 - h(\mu_{-(j-1), M})) q_{W-(j-1)} \quad (j = 2, \dots, \infty). \quad [38]$$

Eq. (31)-Eq. (38) includes the infinite summations. Because these infinite summations cannot be analytically calculated, one should set a cutoff of the summations in the numerical calculation of Eq. (31)-Eq. (38).

Fig. 2-B in the main manuscript shows an example of $\Phi_W^*(p_{WW}, p_{WM})$. In this example, the elements in $\{(\mu_{j, W}, \mu_{j, M})\}_{j \in \mathbb{Z} \setminus \{0\}}$ are all different for different j , and thus the labeling in this study is at least necessary for the description of the reputation structure. This figure also shows that the obtained analytical solutions well approximate simulations of the image matrix.

3. Calculation of expected payoff

In order to consider an evolutionary process, we derive expected payoffs of wild-types W and mutants M from joint probability distribution of goodnesses, i.e., Φ_W^* and Φ_M^* . In the limit that mutants are rare $\delta \rightarrow 0$, the expected payoffs of the wild-types

70 u_W and mutants u_M are given by

$$\begin{aligned} u_W &= (b - c)\bar{p}_{WW}, \\ u_M &= b\bar{p}_{MW} - c\bar{p}_{WM}, \end{aligned} \quad [39]$$

72 Here, $\bar{p}_{AA'}$ indicates the average goodnesses of $p_{AA'}$, i.e., described as

$$\begin{aligned} \bar{p}_{WW} &= \int_0^1 \int_0^1 p_{WW} \Phi_W^*(p_{WW}, p_{WM}) dp_{WW} dp_{WM}, \\ \bar{p}_{WM} &= \int_0^1 \int_0^1 p_{WM} \Phi_W^*(p_{WW}, p_{WM}) dp_{WW} dp_{WM}, \\ \bar{p}_{MW} &= \int_0^1 \int_0^1 p_{MW} \Phi_M^*(p_{MW}, p_{MM}) dp_{MW} dp_{MM}, \end{aligned} \quad [40]$$

74 This average goodness can be analytically calculated by Gaussian approximation of $\Phi_A^*(p_{AW}, p_{AM})$. According to the conditions
75 for the ESS, mutants can invade the population of wild-types if $u_W > u_M$.

76 Regions, where a mutant norm can invade a wild-type norm, are given by Fig. S2. From this figure, we can obtain the
77 invasibilities for each b/c , as shown in Fig. 3-A in the main manuscript.

		mutant												
		01	02	03	04	05	07	08	09	12	13	14	15	16
wild-type	01		always	always	always	always	always	always	always	always	always	always	always	always
	02	> 5.964		< 6.811	< 6.434	< 13.13	< 11.58	< 10.99	< 12.92	< 10.63	< 12.71	< 11.58	< 12.28	< 11.44
	03	> 4.376	> 6.352		< 1.920	never	< 1.306	< 1.441	< 1.134	< 1.457	< 1.261	< 1.402	< 1.306	< 1.424
	04	> 1.250	> 1.250	> 1.250		> 1.250		< 1.250	> 1.250	< 1.250		< 1.250	< 1.250	< 1.250
	05	never	never	never	< 2.789		always	always	always	always	always	always	always	always
	07	never	never	never		never		always	never	always		always	always	always
	08	> 11.44	> 10.99	> 10.63	> 6.434	> 11.58	> 11.58		> 12.28	> 6.811	> 12.71	> 13.13	> 12.92	< 5.964
	09	never	never	always	always	never	always	always		always	always	always	always	always
	12	> 1.424	> 1.441	> 1.457	> 1.920	> 1.402	> 1.306	< 6.352	> 1.306		> 1.261	always	> 1.134	< 4.376
	13	never	never	never		never		always	never	always		always	always	always
	14	never	never	never	> 2.789	never	never	always	never	always	never		never	always
	15	never	never	never	never	never	never	always	never	never	never	always		always
	16	never	never	never	never	never	never	never	never	never	never	never	never	

Fig. S2. Regions for possible invasions in the evolutionary processes. The row and column indicate the wild-type and mutant norms, respectively. The matrix shows the region of $b/c (> 1)$ where the mutant invades the wild-type. In some pairs of wild-type and mutant norms, the mutant always or never succeeds in invading the wild-types for all $b/c (> 1)$. The calculation is based on $e_1 = 0$ and $e_2 = 0.1$.

78 4. Calculation of equilibrium state in public reputation

79 In this section, we derive the equilibrium distribution of reputations under public assessment, based on the previous study (1).
80 The basic setting is the same whether the reputation is publicly shared or privately held. We assume a population of size N
81 which consists of mutants with norm M and wild-type individuals with norm $W \neq M$. A donor and a recipient are randomly
82 chosen every round. The donor chooses cooperation to the good recipient and defection to the bad recipient. Here, the donor
83 erroneously chooses the opposite action to the intended one with probability $0 \leq e_1 < 1/2$. Then, all the individuals update
84 their reputations of the donor. The difference between the public and private reputation cases is seen in the observers' ways
85 to update reputations. We assume that one mutant observer and one wild-type observer are chosen as representatives of
86 each norm, and each gives a good or bad reputation to the donor according to its norm. Here, each representative observer
87 commits an assignment error independently, in which case it erroneously assigns the opposite reputation to the intended one
88 with probability $0 < e_2 < 1/2$. (such an assessment error was not assumed in (1)) Then, all the individuals with the same
89 norm copy the reputation of the donor assigned by their representative. Thus, the reputation of the same individual, even an
90 erroneously assigned one, is shared among all the individuals with the same norm. In other words, each individual at any given
91 time has two reputations, one is shared by all the mutant individuals, and the other is shared by all the wild-type individuals
92 in the population.

93 Here we specifically consider the situation where rare mutants with norm $M = S_{16}(\text{ALLB})$ invades a wild-type population
 94 with norm $W \neq M$. We use the same definition of $p_{AA'}$, i.e., goodness of an individual with norm A in the eyes of norm A'
 95 users. Because reputations are public, $p_{AA'}$ can be either 1 (the individual is assigned as good from all) or 0 (the individual is
 96 assigned as bad from all). Below we will derive $\bar{p}_{AA'}$, the probability that a norm A user has a good reputation in the eyes of
 97 norm A' users.

98 Since the mutant norm is ALLB, the probability that mutants assign a good reputation to the donor is always $a_M^{\text{GC}} = a_M^{\text{BC}} =$
 99 $a_M^{\text{GD}} = a_M^{\text{BD}} = e_2$. Thus we obtain $\bar{p}_{WM} = \bar{p}_{MM} = e_2$.

100 Next we aim to solve the equilibrium average goodnesses in the eyes of wild-types, \bar{p}_{WW} and \bar{p}_{MW} . First, let us calculate
 101 \bar{p}_{WW} , which is relevant when the donor and the observer use norm W . Note that we can assume that the recipient uses norm
 102 W , because mutants are rare. \bar{p}_{WW} should satisfy

$$103 \quad \bar{p}_{WW} = \bar{p}_{WW} \{ (1 - e_1) a_W^{\text{GC}} + e_1 a_W^{\text{GD}} \} + (1 - \bar{p}_{WW}) \{ e_1 a_W^{\text{BC}} + (1 - e_1) a_W^{\text{BD}} \}, \quad [41]$$

104 The equality between the left- and right-hand sides shows that the proportion of good individuals balances before and after
 105 updating the chosen donor's reputation. In the right-hand side, \bar{p}_{WW} and $1 - \bar{p}_{WW}$ in the first and the second terms indicate
 106 the probabilities that a randomly chosen recipient of norm W is good or bad from the viewpoint of norm W , respectively.
 107 When the recipient is good, the donor chooses cooperation or defection with probabilities $(1 - e_1)$ and e_1 . Then, a_W^{GC} and a_W^{GD}
 108 indicate the probabilities that the cooperating or defecting donor receives a good reputation from observers of norm W . When
 109 the recipient is bad, the donor chooses cooperation or defection with probabilities e_1 and $(1 - e_1)$. Then, a_W^{BC} and a_W^{BD} indicate
 110 the probabilities that the cooperating or defecting donor receives a good reputation from observers of norm W . The solution is

$$111 \quad \bar{p}_{WW} = \frac{(1 - e_1) a_W^{\text{BD}} + e_1 a_W^{\text{BC}}}{1 - \{ (1 - e_1) (a_W^{\text{GC}} - a_W^{\text{BD}}) + e_1 (a_W^{\text{GD}} - a_W^{\text{BC}}) \}}. \quad [42]$$

Second, let us calculate \bar{p}_{MW} , which is relevant when the donor uses norm M and the observer uses norm W . Note that we
 can once again assume that the recipient uses norm W because mutants are rare. \bar{p}_{MW} should satisfy

$$112 \quad \bar{p}_{MW} = \bar{p}_{WW} \{ h(e_2) a_W^{\text{GC}} + (1 - h(e_2)) a_W^{\text{GD}} \} + (1 - \bar{p}_{WW}) \{ h(e_2) a_W^{\text{BC}} + (1 - h(e_2)) a_W^{\text{BD}} \}. \quad [43]$$

113 Here, \bar{p}_{WW} and $(1 - \bar{p}_{WW})$ in the first and second terms of the right-hand side indicate the probabilities that the recipient is
 114 good or bad from the viewpoint of norm W , respectively. In both terms, $h(e_2) (= e_2(1 - e_1) + (1 - e_2)e_1)$ and $1 - h(e_2)$ are the
 115 probabilities that the donor with norm M executes cooperation or defection, which is independent of whether the recipient is
 116 good or bad from the viewpoint of norm W . In the first term, a_W^{GC} and a_W^{GD} indicate the probabilities that the cooperating
 117 and defecting donor receives a good reputation from the observers of norm W . In the second term, a_W^{BC} and a_W^{BD} indicate the
 118 probabilities that the cooperating and defecting donor receives a good reputation from the observers of norm W .

We summarize the solutions, \bar{p}_{WW} and \bar{p}_{MW} , in Table S2.

S_k	\bar{p}_{WW}		\bar{p}_{MW}	$\bar{p}_{WW} _{e_1=0}$	$\bar{p}_{MW} _{e_1=0}$
S_{01}	$1 - e_2$	$=$	$1 - e_2$	$1 - e_2$	$1 - e_2$
S_{02}	$\frac{e_1 + e_2 - 2e_1e_2}{e_1 + 2e_2 - 2e_1e_2}$	$<$	$\frac{e_1 + e_2 - 2e_1e_2 + e_2^2 - 2e_1e_2^2 - 2e_2^3 + 4e_1e_2^3}{e_1 + 2e_2 - 2e_1e_2}$	$\frac{1}{2}$	$\frac{1}{2} + \frac{1}{2}e_2 - e_2^2$
S_{03}	$\frac{1 - e_2}{1 + e_1 - 2e_1e_2}$	$>$	$\frac{(1 - e_2)(2e_1 + 3e_2 - 6e_1e_2 - 2e_2^2 + 4e_1e_2^2)}{1 + e_1 - 2e_1e_2}$	$1 - e_2$	$3e_2 - 5e_2^2 + 2e_2^3$
S_{04}	$\frac{1}{2}$	$>$	$\frac{e_1 + 2e_2 - 4e_1e_2 - 2e_2^2 + 4e_1e_2^2}{1 - e_1 - e_2 + 2e_1e_2}$	$\frac{1}{2}$	$2e_2 - 2e_2^2$
S_{05}	$\frac{1 - e_1 - e_2 + 2e_1e_2}{1 - e_1 + 2e_1e_2}$	$>$	$\frac{1 - e_1 - e_2 + 2e_1e_2 - e_2^2 + 2e_1e_2^2 + 2e_2^3 - 4e_1e_2^3}{1 - e_1 + 2e_1e_2}$	$1 - e_2$	$1 - e_2 - e_2^2 + 2e_2^3$
S_{06}	$\frac{1}{2}$	$=$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$
S_{07}	$1 - e_1 - e_2 + 2e_1e_2$	$>$	$2e_1 + 3e_2 - 2e_1^2 - 12e_1e_2 - 6e_2^2 + 12e_1^2e_2 + 24e_1e_2^2 + 4e_2^3 - 24e_1^2e_2^2 - 16e_1e_2^3 + 16e_1^2e_2^3$	$1 - e_2$	$3e_2 - 6e_2^2 + 4e_2^3$
S_{08}	$\frac{e_2}{e_1 + 2e_2 - 2e_1e_2}$	$>$	$\frac{e_2(2e_1 + 3e_2 - 6e_1e_2 - 2e_2^2 + 4e_1e_2^2)}{e_1 + 2e_2 - 2e_1e_2}$	$\frac{1}{2}$	$\frac{3}{2}e_2 - e_2^2$
S_{09}	$\frac{1 - e_2}{2 - e_1 - 2e_2 + 2e_1e_2}$	$<$	$\frac{(1 - e_2)(2 - 2e_1 - 3e_2 + 6e_1e_2 + 2e_2^2 - 4e_1e_2^2)}{2 - e_1 - 2e_2 + 2e_1e_2}$	$\frac{1}{2}$	$1 - \frac{3}{2}e_2 + e_2^2$
S_{10}	$e_1 + e_2 - 2e_1e_2$	$<$	$2e_1 + 3e_2 - 2e_1^2 - 12e_1e_2 - 6e_2^2 + 12e_1^2e_2 + 24e_1e_2^2 + 4e_2^3 - 24e_1^2e_2^2 - 16e_1e_2^3 + 16e_1^2e_2^3$	e_2	$3e_2 - 6e_2^2 + 4e_2^3$
S_{11}	$\frac{1}{2}$	$=$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$
S_{12}	$\frac{e_1 + e_2 - 2e_1e_2}{1 + e_1 - 2e_1e_2}$	$<$	$\frac{e_1 + 2e_2 - 4e_1e_2 - 3e_2^2 + 6e_1e_2^2 + 2e_2^3 - 4e_1e_2^3}{1 + e_1 - 2e_1e_2}$	e_2	$2e_2 - 3e_2^2 + 2e_2^3$
S_{13}	$\frac{1}{2}$	$<$	$\frac{1 - e_1 - 2e_2 + 4e_1e_2 + 2e_2^2 - 4e_1e_2^2}{1 - e_1 - 2e_2 + 4e_1e_2}$	$\frac{1}{2}$	$1 - 2e_2 + 2e_2^2$
S_{14}	$\frac{e_2}{1 - e_1 + 2e_1e_2}$	$<$	$\frac{e_2(2 - 2e_1 - 3e_2 + 6e_1e_2 + 2e_2^2 - 4e_1e_2^2)}{1 - e_1 + 2e_1e_2}$	e_2	$2e_2 - 3e_2^2 + 2e_2^3$
S_{15}	$\frac{1 - e_1 - e_2 + 2e_1e_2}{2 - e_1 - 2e_2 + 2e_1e_2}$	$>$	$\frac{1 - e_1 - e_2 + 2e_1e_2 + 3e_2^2 - 6e_1e_2^2 - 2e_2^3 + 4e_1e_2^3}{2 - e_1 - 2e_2 + 2e_1e_2}$	$\frac{1}{2}$	$\frac{1}{2} - \frac{1}{2}e_2 + e_2^2$
S_{16}	e_2	$=$	e_2	e_2	e_2

Table S2. Analytical solution of reputation structure under public assessment. Here, the equality (i.e., =) and inequality (i.e., > or <) signs show the relations between \bar{p}_{WW} and \bar{p}_{MW} for all of $0 \leq e_1 < 1/2$ and $0 < e_2 < 1/2$.

119 Based on Table S2, we can see how the reputation structure differs between the public and private reputation cases. The
 120 reputation structure for norms S_{03} (SS) and S_{07} (SJ) are illustrated in Fig. 3-D and E in the main manuscript, while that of
 121 S_{08} (SH) is in Fig. S3.

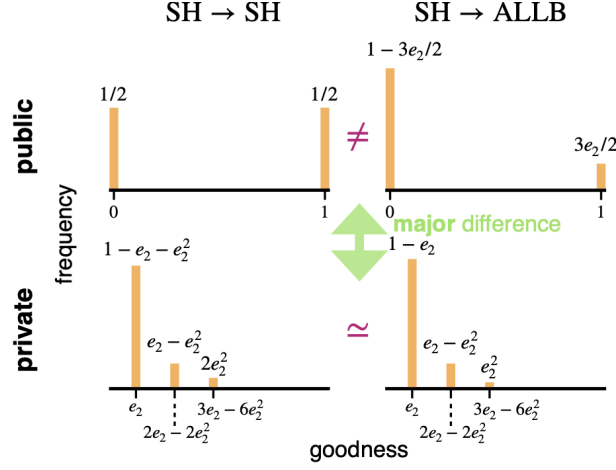


Fig. S3. Illustration of how the wild-type SH gives reputations to the self and mutant ALLB norms. In each panel, the horizontal and vertical axes indicate the goodness and its frequency, respectively. Positions and heights of bars are correct only up to order e_2^2 . By comparing the upper panels with the lower ones, we can see that the reputation from SH differs significantly between the public and private reputation cases. In the private reputation case, SH still manages to distinguish the self norm with ALLB, but only with the difference of the order of e_2^2 .

122 5. Numerical algorithm and error estimate

123 This section provides how to computationally calculate Eq. (31)-Eq. (34) and Eq. (35)-Eq. (38) with sufficient accuracy.

Instead of Eq. (31)-Eq. (34), we aim to compute

$$Q_{W+1} := 1, \quad [44]$$

$$Q_{W+j} := h(\mu_{+(j-1),W})Q_{W+(j-1)} \quad (j = 2, \dots, \infty), \quad [45]$$

$$Q_{W-1} := \sum_{j=1}^{\infty} (1 - h(\mu_{+j,W}))Q_{W+j}, \quad [46]$$

$$Q_{W-j} := (1 - h(\mu_{-(j-1),W}))Q_{W-(j-1)} \quad (j = 2, \dots, \infty), \quad [47]$$

(see Fig. S4 for the illustration of this computation). Via these equations, we obtain q_{Wj} by rescaling Q_{Wj} as

$$q_{Wj} = \frac{Q_{Wj}}{\sum_{k=\pm 1}^{\pm \infty} Q_{Wk}}, \quad [48]$$

which satisfies Eq. (31)-Eq. (34). We should also obtain average goodnesses

$$\bar{p}_{WA} = \frac{\sum_{j=\pm 1}^{\pm \infty} Q_{Wj} \mu_{j,A}}{\sum_{j=\pm 1}^{\pm \infty} Q_{Wj}}, \quad [49]$$

124 in order to obtain Fig. S2.

In a practical computer simulation, we approximate Eq. (44)-Eq. (47) by

$$\hat{Q}_{W+1} := 1, \quad [50]$$

$$\hat{Q}_{W+j} := h(\mu_{+(j-1),W})\hat{Q}_{W+(j-1)} \quad (j = 2, \dots, j_{\max}), \quad [51]$$

$$\hat{Q}_{W+j} := 0 \quad (j = j_{\max} + 1, \dots, \infty), \quad [52]$$

$$\hat{Q}_{W-1} := \sum_{j=1}^{\infty} (1 - h(\mu_{+j,W}))\hat{Q}_{W+j} = \sum_{j=1}^{j_{\max}} (1 - h(\mu_{+j,W}))\hat{Q}_{W+j}, \quad [53]$$

$$\hat{Q}_{W-j} := (1 - h(\mu_{-(j-1),W}))\hat{Q}_{W-(j-1)} \quad (j = 2, \dots, j_{\max}), \quad [54]$$

$$\hat{Q}_{W-j} := 0 \quad (j = j_{\max} + 1, \dots, \infty), \quad [55]$$

with sufficient large $j_{\max} (= 10^4)$ (see Fig. S4 for the illustration of this computation). We will show below that these computationally obtained \hat{Q}_{W_j} well approximate Q_{W_j} . Note that in the following calculations we use the fact that

$$e_2 \leq \mu_{j,A} \leq 1 - e_2 \quad [56]$$

125 holds for all $j = \pm 1, \dots, \pm\infty$ and A .

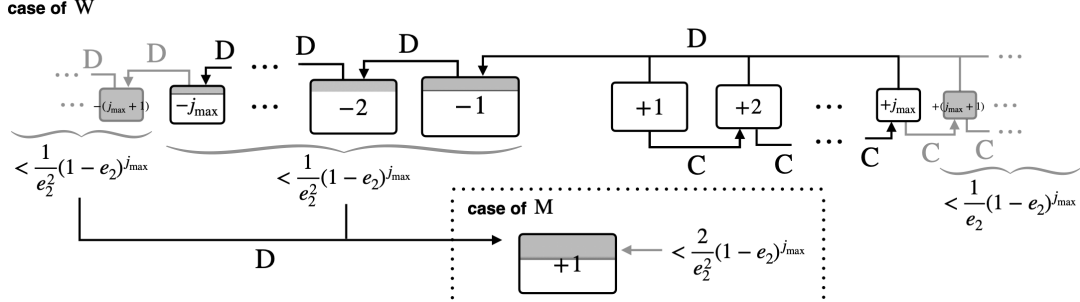


Fig. S4. An illustration of numerical algorithm and error estimation of the masses. The black and gray arrows show how theoretical calculations of Eq. (44)-Eq. (47) are performed, whereas only black arrows are relevant in the computation of Eq. (50)-Eq. (55). Each box shows the size of Q_{W_j} . The gray part in each box shows the size of approximation error, $Q_{W_j} - \hat{Q}_{W_j}$. Apart from Q_{W_j} , the area surrounded by dots show the calculation of Q_{M+1} .

From the definition, we obtain

$$Q_{W+j} - \hat{Q}_{W+j} = 0 \quad (j = 1, \dots, j_{\max}). \quad [57]$$

Then, we obtain

$$Q_{W+j} = Q_{W+1} \prod_{k=1}^{j-1} \mu_{+k,W} \leq (1 - e_2)^{j-1}, \quad [58]$$

$$\Rightarrow \sum_{j=j_{\max}+1}^{\infty} (Q_{W+j} - \hat{Q}_{W+j}) = \sum_{j=j_{\max}+1}^{\infty} Q_{W+j} \leq \sum_{j=j_{\max}+1}^{\infty} (1 - e_2)^{j-1} = \frac{1}{e_2} (1 - e_2)^{j_{\max}}. \quad [59]$$

Then, we have

$$Q_{W-1} = \sum_{j=1}^{\infty} Q_{W+j} (1 - \mu_{+j,W}) = \underbrace{\sum_{j=1}^{j_{\max}} Q_{W+j} (1 - \mu_{+j,W})}_{=\hat{Q}_{W-1}} + \sum_{j=j_{\max}+1}^{\infty} Q_{W+j} (1 - \mu_{+j,W}) \leq \hat{Q}_{W-1} + \frac{1}{e_2} (1 - e_2)^{j_{\max}}, \quad [60]$$

$$\Rightarrow Q_{W-j} = Q_{W-1} \prod_{k=1}^{j-1} (1 - \mu_{-k,W}) = \underbrace{\hat{Q}_{W-1} \prod_{k=1}^{j-1} (1 - \mu_{-k,W})}_{=\hat{Q}_{W-j}} + (Q_{W-1} - \hat{Q}_{W-1}) \prod_{k=1}^{j-1} (1 - \mu_{-k,W}) \leq \hat{Q}_{W-j} + \frac{1}{e_2} (1 - e_2)^{j_{\max}+j-1}, \quad [61]$$

$$\Rightarrow \sum_{j=1}^{j_{\max}} (Q_{W-j} - \hat{Q}_{W-j}) \leq \frac{1}{e_2^2} (1 - e_2)^{j_{\max}}, \quad [62]$$

We also obtain

$$Q_{W-1} = \sum_{j=1}^{\infty} Q_{W+j} (1 - \mu_{+j,W}) \leq \sum_{j=1}^{\infty} Q_{W+j} \leq \frac{1}{e_2}, \quad [63]$$

$$\Rightarrow Q_{W-j} = Q_{W-1} \prod_{k=1}^{j-1} (1 - \mu_{-k,W}) \leq \frac{1}{e_2} (1 - e_2)^{j-1}, \quad [64]$$

$$\Rightarrow \sum_{j=j_{\max}+1}^{\infty} (Q_{W-j} - \hat{Q}_{W-j}) = \sum_{j=j_{\max}+1}^{\infty} Q_{W-j} \leq \sum_{j=j_{\max}+1}^{\infty} \frac{1}{e_2} (1 - e_2)^{j-1} = \frac{1}{e_2^2} (1 - e_2)^{j_{\max}}. \quad [65]$$

From the above error estimations, we can obtain upper and lower bounds of Eq. (48) and Eq. (49) as

$$\frac{\hat{Q}_{Wj}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} + \frac{3}{e_2^2}(1-e_2)^{j_{\max}}} \leq q_{Wj} \leq \frac{\hat{Q}_{Wj} + \frac{1}{e_2}(1-e_2)^{j_{\max}}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj}}. \quad [66]$$

$$\frac{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} \mu_{j,A}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} + \frac{3}{e_2^2}(1-e_2)^{j_{\max}}} \leq \bar{p}_{WA} \leq \frac{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} \mu_{j,A} + \frac{3}{e_2^2}(1-e_2)^{j_{\max}}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj}}. \quad [67]$$

Here, we used

$$\hat{Q}_{Wj} \leq Q_{Wj} = \hat{Q}_{Wj} + (Q_{Wj} - \hat{Q}_{Wj}) \quad [68]$$

$$\begin{aligned} &\leq \hat{Q}_{Wj} + \underbrace{\max_j (Q_{Wj} - \hat{Q}_{Wj})}_{=Q_{W-1} - \hat{Q}_{W-1} \leq \frac{1}{e_2}(1-e_2)^{j_{\max}}} \leq \hat{Q}_{Wj} + \frac{1}{e_2}(1-e_2)^{j_{\max}}, \end{aligned} \quad [69]$$

$$\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} \leq \sum_{j=\pm 1}^{\pm \infty} Q_{Wj} = \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} + \sum_{j=\pm 1}^{\pm \infty} (Q_{Wj} - \hat{Q}_{Wj}) \quad [70]$$

$$\begin{aligned} &= \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} + \underbrace{\sum_{j=1}^{j_{\max}+1} (Q_{Wj} - \hat{Q}_{W-j})}_{\leq \frac{1}{e_2}(1-e_2)^{j_{\max}}} + \underbrace{\sum_{j=j_{\max}+1}^{\infty} (Q_{W+j} - \hat{Q}_{W+j})}_{\leq \frac{1}{e_2}(1-e_2)^{j_{\max}}} + \underbrace{\sum_{j=j_{\max}+1}^{\infty} (Q_{W-j} - \hat{Q}_{W-j})}_{\leq \frac{1}{e_2}(1-e_2)^{j_{\max}}} \end{aligned} \quad [71]$$

$$\leq \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} + \frac{3}{e_2^2}(1-e_2)^{j_{\max}}, \quad [72]$$

$$\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} \mu_{j,A} \leq \sum_{j=\pm 1}^{\pm \infty} Q_{Wj} \mu_{j,A} \leq \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} \mu_{j,A} + \sum_{j=\pm 1}^{\pm \infty} (Q_{Wj} - \hat{Q}_{Wj}) \quad [73]$$

$$\leq \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Wj} \mu_{j,A} + \frac{3}{e_2^2}(1-e_2)^{j_{\max}}. \quad [74]$$

In the same way, let us consider Eq. (35)-Eq. (38) and compute

$$Q_{M+1} := \sum_{j=1}^{\infty} h(\mu_{-j,M}) Q_{W-j}, \quad [75]$$

$$Q_{M+j} := h(\mu_{+(j-1),M}) Q_{W+(j-1)} \quad (j = 2, \dots, \infty), \quad [76]$$

$$Q_{M-1} := \sum_{j=1}^{\infty} (1 - h(\mu_{+j,M})) Q_{W+j}, \quad [77]$$

$$Q_{M-j} := (1 - h(\mu_{-(j-1),M})) Q_{W-(j-1)} \quad (j = 2, \dots, \infty). \quad [78]$$

Via these equations, we obtain q_{Wj} by rescaling Q_{Mj} as

$$q_{Mj} = \frac{Q_{Mj}}{\sum_{k=\pm 1}^{\pm \infty} Q_{Mk}}, \quad [79]$$

which satisfies Eq. (35)-Eq. (38). We also need to obtain average goodnesses

$$\bar{p}_{MA} = \frac{\sum_{j=\pm 1}^{\pm \infty} Q_{Mj} \mu_{j,A}}{\sum_{j=\pm 1}^{\pm \infty} Q_{Mj}}, \quad [80]$$

In a practical computer simulation, we approximate Eq. (75)-Eq. (78) by

$$\hat{Q}_{M+1} := \sum_{j=1}^{\infty} h(\mu_{-j,M}) \hat{Q}_{W-j} = \sum_{j=1}^{j_{\max}} h(\mu_{-j,M}) \hat{Q}_{W-j}, \quad [81]$$

$$\hat{Q}_{M+j} := h(\mu_{+(j-1),M}) \hat{Q}_{W+(j-1)} \quad (j = 2, \dots, j_{\max}), \quad [82]$$

$$\hat{Q}_{M+j} := 0 \quad (j = j_{\max} + 1, \dots, \infty), \quad [83]$$

$$\hat{Q}_{M-1} := \sum_{j=1}^{\infty} (1 - h(\mu_{+j,M})) \hat{Q}_{W+j} = \sum_{j=1}^{j_{\max}} (1 - h(\mu_{+j,M})) \hat{Q}_{W+j}, \quad [84]$$

$$\hat{Q}_{M-j} := (1 - h(\mu_{-(j-1),M})) \hat{Q}_{W-(j-1)} \quad (j = 2, \dots, j_{\max}), \quad [85]$$

$$\hat{Q}_{M-j} := 0 \quad (j = j_{\max} + 1, \dots, \infty), \quad [86]$$

127 with sufficient large $j_{\max}(= 10^4)$.

By exactly similar calculations, we obtain

$$Q_{M+j} - \hat{Q}_{M+j} = 0 \quad (j = 2, \dots, j_{\max}), \quad [87]$$

$$\sum_{j=j_{\max}+1}^{\infty} (Q_{M+j} - \hat{Q}_{M+j}) \leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}, \quad [88]$$

$$\sum_{j=1}^{j_{\max}} (Q_{M-j} - \hat{Q}_{M-j}) \leq \frac{1}{e_2^2} (1 - e_2)^{j_{\max}}, \quad [89]$$

$$\sum_{j=j_{\max}+1}^{\infty} (Q_{M-j} - \hat{Q}_{M-j}) \leq \frac{1}{e_2^2} (1 - e_2)^{j_{\max}}. \quad [90]$$

The difference from Q_{Wj} exists only in $j = +1$, as

$$Q_{M+1} = \sum_{j=1}^{\infty} Q_{W-j} \mu_{-j,M} = \underbrace{\sum_{j=1}^{j_{\max}} \hat{Q}_{W-j} \mu_{-j,M}}_{=\hat{q}_{M+1}} + \underbrace{\sum_{j=1}^{j_{\max}} (Q_{W-j} - \hat{Q}_{W-j}) \mu_{-j,M}}_{\leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}} + \underbrace{\sum_{j=j_{\max}+1}^{\infty} \hat{Q}_{W-j} \mu_{-j,M}}_{\leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}} \quad [91]$$

$$\leq \hat{Q}_{M+1} + \frac{2}{e_2} (1 - e_2)^{j_{\max}}, \quad [92]$$

128 (see the area surrounded by dots in Fig. S4 for the illustration of this computation).

From the above error estimations, we can obtain upper and lower bounds of Eq. (79) and Eq. (80) as

$$\frac{\hat{Q}_{Mj}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} + \frac{5}{e_2} (1 - e_2)^{j_{\max}}} \leq q_{Mj} \leq \frac{\hat{Q}_{Mj} + \frac{2}{e_2} (1 - e_2)^{j_{\max}}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj}}, \quad [93]$$

$$\frac{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} \mu_{j,A}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} + \frac{5}{e_2} (1 - e_2)^{j_{\max}}} \leq \bar{p}_{MA} \leq \frac{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} \mu_{j,A} + \frac{5}{e_2} (1 - e_2)^{j_{\max}}}{\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj}}. \quad [94]$$

Here, we used

$$\hat{Q}_{Mj} \leq Q_{Mj} = \hat{Q}_{Mj} + (Q_{Mj} - \hat{Q}_{Mj}) \quad [95]$$

$$\leq \hat{Q}_{Mj} + \underbrace{\max_j (Q_{Mj} - \hat{Q}_{Mj})}_{=Q_{M+1} - \hat{Q}_{M+1} \leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}} \leq \hat{Q}_{Mj} + \frac{1}{e_2} (1 - e_2)^{j_{\max}}, \quad [96]$$

$$\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} \leq \sum_{j=\pm 1}^{\pm \infty} Q_{Mj} = \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} + \underbrace{(Q_{M+1} - \hat{Q}_{M+1})}_{\leq \frac{2}{e_2} (1 - e_2)^{j_{\max}}} + \underbrace{\sum_{j=1}^{j_{\max}} (Q_{M-j} - \hat{Q}_{M-j})}_{\leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}} \quad [97]$$

$$+ \underbrace{\sum_{j=j_{\max}+1}^{\infty} (Q_{M+j} - \hat{Q}_{M+j})}_{\leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}} + \underbrace{\sum_{j=j_{\max}+1}^{\infty} (Q_{M-j} - \hat{Q}_{M-j})}_{\leq \frac{1}{e_2} (1 - e_2)^{j_{\max}}} \quad [98]$$

$$\leq \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} + \frac{5}{e_2} (1 - e_2)^{j_{\max}}, \quad [99]$$

$$\sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} \mu_{j,A} \leq \sum_{j=\pm 1}^{\pm \infty} Q_{Mj} \mu_{j,A} \leq \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} \mu_{j,A} + \sum_{j=\pm 1}^{\pm \infty} (Q_{Mj} - \hat{Q}_{Mj}) \quad [100]$$

$$\leq \sum_{j=\pm 1}^{\pm j_{\max}} \hat{Q}_{Mj} \mu_{j,A} + \frac{5}{e_2} (1 - e_2)^{j_{\max}}. \quad [101]$$

References

1. S Uchida, K Sigmund, The competition of assessment rules for indirect reciprocity. *J. Theor. Biol.* **263**, 13–19 (2010).