



International Institute for
Applied Systems Analysis
Schlossplatz 1
A-2361 Laxenburg, Austria

Tel: +43 2236 807 342
Fax: +43 2236 71313
E-mail: publications@iiasa.ac.at
Web: www.iiasa.ac.at

Interim Report

IR-12-067

The evolution of cooperation by social exclusion

Tatsuya Sasaki (sasakit@iiasa.ac.at)
Satoshi Uchida

Approved by

Ulf Dieckmann
Director, Evolution and Ecology Program

February 2015

Interim Reports on work of the International Institute for Applied Systems Analysis receive only limited review. Views or opinions expressed herein do not necessarily represent those of the Institute, its National Member Organizations, or other organizations supporting the work.

The evolution of cooperation by social exclusion

Sasaki T, Uchida S. 2013 The evolution of cooperation by social exclusion. *Proc. R. Soc. B.* **280**, 20122498; doi:10.1098/rspb.2012.2498 (published 5 December 2012). Reprint is available from <http://rspb.royalsocietypublishing.org/cgi/reprint/rspb.2012.2498?ijkey=X87Vy0Xhm7hwMqH&keytype=ref>

Tatsuya Sasaki^{a,1} and Satoshi Uchida^b

^aEvolution and Ecology Program, International Institute for Applied Systems Analysis,
Schlossplatz 1, 2631 Laxenburg, Austria

^bResearch Center, RINRI Institute, Misaki-cho 3-1-10, Chiyoda-ku, 101-8385 Tokyo, Japan

¹To whom correspondence should be addressed.

E-mail: sasakit@iiasa.ac.at

6 December 2012

Preprint version 2.0

14 **Summary:** The exclusion of freeriders from common privileges or public acceptance is widely
15 found in the real world. Current models on the evolution of cooperation with incentives mostly
16 assume peer sanctioning, whereby a punisher imposes penalties on freeriders at a cost to itself.
17 It is well known that such costly punishment has two substantial difficulties. First, a rare
18 punishing cooperator barely subverts the asocial society of freeriders, and second, natural
19 selection often eliminates punishing cooperators in the presence of non-punishing cooperators
20 (namely, “second-order” freeriders). We present a game-theoretical model of social exclusion
21 in which a punishing cooperator can exclude freeriders from benefit sharing. We show that
22 such social exclusion can overcome the above-mentioned difficulties even if it is costly and
23 stochastic. The results do not require a genetic relationship, repeated interaction, reputation, or
24 group selection. Instead, only a limited number of freeriders are required to prevent the second-
25 order freeriders from eroding the social immune system.

26 **Key words:** evolution of cooperation; ostracism; costly punishment; second-order freerider;
27 public goods; evolutionary game theory

28

29 **1. Introduction**

30 We frequently engage in voluntary joint enterprises with nonrelatives, activities that are
31 fundamental to society. The evolution of cooperative behaviors is an important issue because
32 without any supporting mechanism [1], natural selection often favours those that contribute
33 less at the expense of those that contribute more. A minimal situation could easily cause the
34 ruin of a commune of cooperators, namely, the “tragedy of the commons” [2]. Here we
35 consider different types of punishment, such as a monetary fine (e.g., [3–7]) and ostracism (e.g.,
36 [8–11]), for the evolution of cooperation. Punishment can reduce the expected payoff for the
37 opponent, and subsequently, change natural selection preferences, to encourage additional
38 contributions to communal efforts [12]. Our model looks at this situation, because “very little
39 work has addressed questions about the form that punishment is likely to take in reality and
40 about the relative efficacy of different types of punishment” [13].

41 Here, we choose to focus on social exclusion, which is a common and powerful tool to penalise
42 deviators in human societies, and includes behaviors such as eviction, shunning and ignoring
43 [14–16]. For self-sustaining human systems, indeed, the ability to distinguish among
44 individuals and clarify who should participate in the sharing of communal benefits is crucial
45 and expected (of its members) [17]. A specific example is found in the case of traffic violators
46 who are punished, often strictly by suspending or revoking their driver license for public roads.
47 Among non-humans, shunning through partner switching is a common mechanism for inequity
48 aversion and cooperation enforcement [13,18,19]. Experimental studies have shown, for
49 instance, that chimpanzees can use a mechanism to exclude less cooperative partners from
50 potential collaborations [20], or that reef fish will terminate interaction with cleaner fish that
51 cheat by eating the host’s mucus rather than parasites [21].

52 In joint enterprises, by excluding freeriders from benefit sharing, the punishers can naturally
53 benefit, because such exclusion often decreases the number of beneficiaries, with little effect
54 on the total benefit. Consider the example of the division of a pie provided by some volunteers

55 to a group. If a person is one of the volunteers, it may be justifiable in terms of fairness to
56 suggest or even force freeriders to refrain from sharing in the pie. Although excluding
57 freeriders can be stressful, it increases the share of the pie for the contributors, including the
58 person who performs the actual exclusion. If the situation calls for it, the excluded freerider's
59 share of the group benefits may separately be redistributed among the remaining members in
60 the group [22,23]. Therefore, in either case, the excluded member will obtain nothing from the
61 joint enterprise and the exclusion causes immediate increases in the payoff for the punisher and
62 also the other remaining members in the group.

63 This is a “self-serving” form of punishment [13,18]. It is of importance that if the cost of
64 excluding is smaller than the reallocated benefit, social exclusion can provide immediate net
65 benefits even to the punisher. This can potentially motivate the group members to contribute to
66 the exclusion of freeriders, however, our understanding of how cooperation unfolds through
67 social exclusion is still “uncharted territory” [24].

68 Most game-theoretical works on cooperation with punishment have focused on other forms of
69 punishment, for example, costly punishment that reduces the payoffs of both the punishers and
70 those who are punished. As is well known, costly punishment poses fundamental puzzles with
71 regard to its emergence and maintenance. First of all, costly punishment is unlikely to emerge
72 in a sea of freeriders, in which almost all freeriders are unaffected, and a rare punisher would
73 have to decrease in its payoff through punishing the left and right [18,25–27]. Moreover,
74 although initially prevalent, punishers can stabilise cooperation, while non-punishing
75 cooperators (so-called “second-order freeriders”) can undermine full cooperation once it is
76 established [3,13,17,24,29].

77 In terms of self-serving punishments, however, we have only started to confront the puzzles
78 that emerge in these scenarios. We ask here, what happens if social exclusion is applied?: that
79 is, do players move toward excluding others?, and can freeriders be eliminated? Or, will others
80 in the group resist? Our main contribution is to provide a detailed comparative analysis for

81 social exclusion and costly punishment, two different types of punishment, from the viewpoint
82 of their emergence and maintenance. With the self-serving function, social exclusion is
83 predicted to more easily emerge and be maintained than costly punishment.

84 Few theoretical works have investigated the conditions under which cooperation can evolve by
85 the exclusion of freeriders. Our model requires no additional modules, such as a genetic
86 relationship, repeated games, reputation, or group selection. Considering these modules is
87 imperative for understanding the evolution of cooperation in realistic settings. In fact, these
88 modules may have already been incorporated in earlier game-theoretical models that included
89 the exclusion of freeriders [30–32], but we are interested in first looking at the most minimal of
90 situations to get at the core relative efficacy of costly punishment versus social exclusion.

91 **2. Game-theoretical model and analysis**

92 To describe these punishment schemes in detail, we begin with standard public good games
93 with a group size of $n \geq 2$ (e.g., [26,33,34]) in an infinitely large, well-mixed population of
94 players. We specifically apply a replicator system [35] for the dynamic analysis, as based on
95 preferentially imitating strategies of the more successful individuals. In the game, each player
96 has two options. The “cooperator” contributes $c > 0$ to a common pool, and the “defector”
97 contributes nothing. The total contribution is multiplied by a factor of $r > 1$ and then shared
98 equally among all (n) group members. A cooperator will thus pay a net cost $\sigma = c(1 - r/n)$
99 through its own contribution. If all cooperate, the group yields the optimal benefit $c(r - 1)$ for
100 each; if all defect, the group does nothing. To adhere to the spirit of the tragedy of the
101 commons, we hereafter assume that $r < n$ holds, in which case a defecting player can improve
102 its payoff by $\sigma > 0$, whatever the coplayers do, and the defectors dominate the cooperators. To
103 observe the robustness for stochastic effects, we also consider an individual-based simulation
104 with a pairwise comparison process [36,37]. See the electronic supplementary material (ESM)
105 for these details.

106 **(a) Costly punishment**

107 We then introduce a third strategy, “punisher”, which contributes c , and moreover, punishes
 108 the defectors. Punishing incurs a cost $\gamma > 0$ per defector to the punisher and imposes a fine
 109 $\beta > 0$ per punisher on the defector. We denote by x , y , and z the frequencies of the cooperator
 110 (C), defector (D), and punisher (P), respectively. Thus, $x, y, z \geq 0$ and $x + y + z = 1$. Given the
 111 expected payoffs P_S for the three strategies ($S = C, D$, and P), the replicator system is written
 112 by

$$113 \quad \dot{x} = x(P_C - \bar{P}), \quad \dot{y} = y(P_D - \bar{P}), \quad \dot{z} = z(P_P - \bar{P}), \quad (2.1)$$

114 where $\bar{P} := xP_C + yP_D + zP_P$ describes the average payoff in the entire population. Three
 115 homogeneous states ($x = 1$, $y = 1$, and $z = 1$) are equilibria. Indeed,

$$116 \quad P_C = \frac{rc}{n}(n-1)(x+z) - \sigma, \quad (2.2a)$$

$$117 \quad P_D = \frac{rc}{n}(n-1)(x+z) - \beta(n-1)z, \quad (2.2b)$$

$$118 \quad P_P = \frac{rc}{n}(n-1)(x+z) - \sigma - \gamma(n-1)y. \quad (2.2c)$$

119 Here the common first term denotes the benefit that resulted from the expected $(n-1)(x+z)$
 120 contributors among the $(n-1)$ coplayers, and $\beta(n-1)z$ and $\gamma(n-1)y$ give the expected fine on
 121 a defector and expected cost to a punisher, respectively.

122 First, consider only the defectors and punishers (figure 1). Thus, $y + z = 1$, and the replicator
 123 system reduces to $\dot{z} = z(1-z)(P_P - P_D)$. Solving $P_P = P_D$ results in that, if the interior
 124 equilibrium R between the two strategies exists, it is uniquely determined by

$$125 \quad z = 1 - \frac{(n-1)\beta - \sigma}{(n-1)(\beta + \gamma)}. \quad (2.3)$$

126 The point R is unstable. If the fine is much smaller: $\beta < \sigma/(n-1) =: \beta_0$, punishment has no
 127 effect on defection dominance, or otherwise, R appears and the dynamics turns into bistable
 128 [33,34]: R separates the state space into basins of attraction of the different homogeneous
 129 states for both the defector and excluder. The smaller γ or larger β , the more the coordinate
 130 of R shifts to the defector end: the more relaxed the initial condition required to establish a
 131 punisher population (figure 1a). Note that a rare punisher is incapable of invading a defector

132 population because the resident defectors, almost all unpunished, earn 0 on average, and the
133 rare punisher does $-\sigma - \gamma(n-1) < 0$.

134 Next, consider all of the cooperators, defectors, and punishers (figure 1b). Without defectors,
135 no punishing cost arises. Thus, no natural selection occurs between the cooperators and
136 punishers, and the edge between the cooperators and punishers ($x + z = 1$) consists of fixed
137 points. A segment consisting of these fixed points with $z > \beta_0/\beta$ is stable against the invasion
138 of rare defectors, and the other segment not so [33,34]. Therefore, this stable segment appears
139 on the edge PC if and only if the edge PD is bistable. We denote by K_0 the boundary point with
140 $z = \beta_0/\beta$. There can thus be two attractors: the vertex D and segment PK_0 . The smaller γ or
141 larger β , the broader the basin of attraction for the mixture states of the contributors. That is,
142 the higher the punishment efficiency, the more relaxed the initial condition required to
143 establish a cooperative state. This may collaborate with evidence from recent public-good
144 experiments [38–40], which suggest the positive effects of increasing the punishment
145 efficiency on average cooperation.

146 However, the stability of PK_0 is not robust for small perturbations of the population. Since
147 $P_p < P_c$ holds in the interior space, an interior trajectory eventually converges to the boundary,
148 and $d(z/x)/dt = (z/x)(P_p - P_c) < 0$: the frequency ratio of the punishers to cooperators
149 decreases over time. Thus, if rare defectors are introduced, for example by mutation or
150 immigration, into a stable population of the two types of contributors, the punishers will
151 gradually decline for each elimination of the defectors. Such small perturbations push the
152 population into an unstable regime around K_0C , where the defectors can invade the population
153 and then take it over. See figure S1 of ESM and also [26] for individual-based simulations.

154 **(b) Social exclusion**

155 We turn next to social exclusion. The third strategy is now replaced with the excluder (E) that
156 contributes c and also tries to exclude defectors from sharing benefits at a cost to itself of
157 $\bar{\gamma} > 0$ per defector. The multiplied contribution is shared equally among the remaining

158 members in the group. We assume that an excluder succeeds in excluding a defector with the
 159 probability $\bar{\beta}$ and that the excluded defector earns nothing. For simplicity, we conservatively
 160 assume that the total sanctioning cost for an excluder is given by $\bar{\gamma}$ times the number of
 161 defectors in a group, whatever others do.

162 We focus on perfect exclusion with $\bar{\beta} = 1$: exclusion never fails. Under this condition, however,
 163 we can analyse the nature of social exclusion considered for cooperation. Indeed, we formalise
 164 the expected payoffs, as follows:

$$165 \quad P_C = c(r-1) - (1-z)^{n-1} \frac{rc}{n} (n-1) \frac{y}{1-z}, \quad (2.4a)$$

$$166 \quad P_D = (1-z)^{n-1} \frac{rc}{n} (n-1) \frac{x}{1-z}, \quad (2.4b)$$

$$167 \quad P_E = c(r-1) - \bar{\gamma}(n-1)y. \quad (2.4c)$$

168 Equation (2.4c) describes that the excluder can constantly receive the group optimum $c(r-1)$
 169 at the exclusion cost expected as $\bar{\gamma}(n-1)y$. In equations (2.4a) and (2.4b), $(1-z)^{n-1}$ denotes the
 170 probability that we find no excluder in the $(n-1)$ coplayers, and if so, $(n-1)y/(1-z)$ and
 171 $(n-1)x/(1-z)$ give the expected numbers of the defectors and cooperators, respectively,
 172 among the coplayers. Hence, the second term of equations (2.4a) specifies an expected benefit
 173 that could have occurred without freeriding, and equation (2.4b) describes an expected amount
 174 that a defector has nibbled from the group benefit, in the group with no excluder. The expected
 175 payoffs for any $\bar{\beta}$ are formalised in ESM.

176 First, the dynamics between the excluders and defectors can only exhibit bi-stability or
 177 excluder dominance for $\bar{\beta} = 1$ (figure 2a). Considering that $P_D = 0$ holds for whatever the
 178 fraction of excluders, solving $P_E = 0$ gives that, if the interior equilibrium R exists, it is
 179 uniquely determined by

$$180 \quad z = 1 - \frac{(r-1)c}{(n-1)\bar{\gamma}}. \quad (2.5)$$

181 The point R is unstable. As before, for larger values of $\bar{\gamma}$, the dynamics between the two
 182 strategies have been bistable. The smaller the value of $\bar{\gamma}$, the larger the basin of attraction to

183 the vertex E. In contrast to costly punishment, an excluder population can evolve, irrespective
 184 of the initial condition, for sufficiently small values of $\bar{\gamma}$. When decreasing $\bar{\gamma}$ beyond a
 185 threshold value, R exits at the vertex D, and thus, the current dynamics of bi-stability turns into
 186 excluder dominance. From substituting $z = 0$ into equation (2.5), the threshold value is
 187 calculated as $\bar{\gamma}_0 = (r-1)c/(n-1)$. We note that the dynamics exhibit defector dominance no
 188 matter what $\bar{\gamma}$, if $\bar{\beta}$ is smaller than z_0 , which is from solving $(1-\bar{\beta})^{n-1}rc(n-1)/n > c(r-1)$:
 189 the unexcluded rare defector is better off than the resident excluders.

190 Next, consider all three strategies (figure 2b). Solving $P_C = P_D$ results in

$$191 \quad z = 1 - \left(\frac{n(r-1)}{r(n-1)} \right)^{\frac{1}{n-1}} =: z_0. \quad (2.6)$$

192 By the assumption $r < n$, we have $0 < z_0 < 1$. Let us denote by K_0 a point at which this line
 193 connects to the edge EC ($x + y = 1$). This edge consists of fixed points, each of which
 194 corresponds to a mixed state of the excluders and cooperators. These fixed points on the
 195 segment EK_0 ($z > z_0$) are stable, and those on the segment K_0C are unstable. Similarly, solving
 196 $P_E = P_C$ gives

$$197 \quad z = 1 - \left(\frac{n\gamma}{rc} \right)^{\frac{1}{n-2}} =: z_1. \quad (2.7)$$

198 We denote by K_1 a point at which the line $z = z_1$ connects to EC. These two lines are parallel,
 199 and thus, there is no generic interior equilibrium.

200 Importantly, the time derivative of z/x is positive in the interior region with $z < z_1$. Therefore,
 201 the dynamics around the segment K_1K_0 are found to be the opposite of costly punishment, if
 202 $z_1 > z_0$ (or otherwise, K_1K_0 has been unstable against rare defectors). In this case, introducing
 203 rare defectors results in that, for each elimination of the defectors, the excluders will gradually
 204 rise along K_1K_0 yet fall along the segment EK_1 . Consequently, with such small perturbations,
 205 the population can remain attracted to the vicinity of K_1 , not converging to D. Moreover, if
 206 $\bar{\gamma} < \bar{\gamma}_0$, the excluders dominate the defectors, and thus, all interior trajectories converge to the

207 segment EK_0 , which appears globally stable (figure 2*b*). This result remains robust for the
208 intermediate exclusion probability (figure 3). See figures S2 and S3 of ESM for individual-
209 based simulations.

210 **3. Discussion**

211 Our results regarding social exclusion show that it can be a powerful incentive and appears in
212 stark contrast to costly punishment. What is the logic behind this outcome? First, it is a fact
213 that the exclusion of defectors can decrease the number of beneficiaries, especially when it
214 does not affect the contributions, thereby increasing the share of the group benefit. Therefore,
215 in a mixed group of excluders and defectors, the excluder's net payoff can become higher than
216 the excluded defector's payoff, which is nothing, especially if the cost to exclude is sufficiently
217 low. If social exclusion is capable of 100% rejection at a cheap cost, it can thus emerge in a sea
218 of defectors and dominate them. In our model, self-serving punishment can emerge even when
219 freeriding is initially prevalent by allowing high net benefits from the self-serving action.

220 Moreover, we find that an increase in the fraction of excluders produces a higher probability of
221 an additional increase in the excluder's payoff. This effect can yield the well-known Simpson's
222 paradox (e.g., [41]): the excluders can obtain a higher average payoff than the cooperators,
223 despite the fact that the cooperators always do better than the excluders for any mixed group of
224 the cooperators, defectors, and excluders. Hence, in the presence of defectors, the replicator
225 dynamics often favour the excluders at the expense of the cooperators. Significantly, if a player
226 may occasionally mutate to a defector, social exclusion is more likely than costly punishment
227 to sustain a cooperative state in which all contribute. In our model, a globally stable,
228 cooperative regime can be sustained when solving the second-order freerider problem by
229 allowing mutation to freeriders.

230 Sanctioning the second-order freeriders has also often been considered for preventing their
231 proliferation [3,29,34,36], although such second-order sanction appears rare in experimental
232 settings [42]. And, allowing for our simple model, it is obvious that in the presence of

233 defectors and cooperators, a second-order punisher that also punishes the cooperators is worse
234 off than the existing punisher, and thus, does not affect defector dominance as in our main
235 model. However, given that excluding more coplayers can cause an additional increase in the
236 share of the group benefit, it is worth exploring whether the second-order excluder that also
237 excludes the cooperators is more powerful than the excluder. Interestingly, our preliminary
238 individual-based investigation often finds that second-order excluders are undermined by the
239 excluders and cooperators, which forms a stable coexistence (figures S4 of EMS): second-
240 order exclusion can be redundant.

241 A fundamental assumption of the model is that defection can be detected with no or little cost.
242 This assumption appears most applicable to local public goods and team production settings in
243 which the coworker's contribution can be easily monitored. However, if the monitoring of co-
244 players for defection imposes a certain cost on the excluders, the cooperators dominate the
245 excluders, and the exclusion-based full cooperation is no longer stable. A typical example is
246 found in a potluck party that will often rotate so that every member takes charge of the party by
247 rotation. This rotation system can promote the equal sharing of the hosting cost; or otherwise,
248 no one would take turn playing host.

249 We assessed by extensive numerical investigations the robustness of our results with respect to
250 the following variants (figures S5 and S6 of EMS). First, we considered a different group size
251 n [3,43], In costly punishment, the stable segment PK_0 expands with n , yet our main results
252 were unaffected: with small perturbations, the population eventually converges to a non-
253 cooperative state in which all freeride. In social exclusion, our results remain qualitatively
254 robust with smaller and larger sizes ($n = 4$ and $n = 10$), but the limit exclusion cost $\bar{\gamma}$ becomes
255 more restricted as n increases. Next, we considered a situation in which a punisher or excluder
256 can choose the number of defectors they sanction. For simplicity, here we assume that each of
257 them sanctions only one [22,44], who is selected randomly from all defectors in the group. Our
258 results remain unaffected, except that social exclusion becomes incapable of emerging in a
259 defector population, in which the payoff of a rare excluder is only given by

260 $rc/(n-1) - c - \bar{\gamma} < 0$. To bring forth the possibility of an emergence, a rare excluder is required
261 to exclude more than $n - rc / (c + \bar{\gamma})$ defectors.

262 Our results spur new questions about earlier studies on the evolution of cooperation with
263 punishment. A fascinating extension is to the social structures through which individuals
264 interact. To date, a large body of work on cooperation has looked at how costly punishment
265 can propagate throughout a social network [45–47]: for example, the interplay of costly
266 punishment and reputation can promote cooperation [48]; strict-and-severe punishment and
267 cooperation can jointly evolve with continuously varying strategies [49]; and evolution can
268 favour anti-social punishment that targets cooperators [50]. Our results show that social
269 exclusion as considered is so simple, yet extremely powerful. That is, even intuitively applying
270 it to previous studies can help us much in understanding how humans and non-humans have
271 been incentivized to exclude freeriders.

272 To resist the exclusion, it is likely that conditional cooperators capable of detecting ostracism
273 (e.g., [8]) evolve. This would then raise the comprehensive cost of exclusion to the excluders
274 because of more difficulties of finding and less opportunities of excluding freeriders. This
275 situation can then result in driving an arms race of the exclusion technique and exclusion
276 detection system. An extensive investigation for understanding joint evolution of these systems
277 is for future work.

278 **Acknowledgements**

279 We thank Karl Sigmund and Voltaire Cang for their comments and suggestions. This study
280 was enabled by financial support by the FWF (Austrian Science Fund) to Ulf Dieckmann at
281 IIASA (TECT I-106 G11), and was also supported by grant RFP-12-21 from the Foundational
282 Questions in Evolutionary Biology Fund.

283

284 **References**

- 285 1. Nowak, M. A. 2012 Evolving cooperation. *J. Theor. Biol.* **299**, 1–8.
286 (doi:10.1016/j.jtbi.2012.01.014)
- 287 2. Hardin, G. 1968 The tragedy of the commons. *Science* **162**, 1243–1248.
288 (doi:10.1126/science.162.3859.1243)
- 289 3. Boyd, R. & Richerson, P. 1992 Punishment allows the evolution of cooperation (or
290 anything else) in sizable groups. *Ethol. Sociobiol.* **13**, 171–195. (doi:10.1016/0162-
291 3095(92)90032-Y)
- 292 4. Fehr, E. & Gächter, S. 2002 Altruistic punishment in humans. *Nature* **415**, 137–140.
293 (doi:10.1038/415137a)
- 294 5. Masclet, D., Noussair, C., Tucker, S. & Villeval, M. -C. 2003 Monetary and
295 nonmonetary punishment in the voluntary contributions mechanism. *Am. Econ. Rev.* **93**,
296 366–380. (doi:10.1257/000282803321455359)
- 297 6. Sigmund, K. 2007 Punish or perish? Retaliation and collaboration among humans.
298 *Trends. Ecol. Evol.* **22**, 593–600. (doi:10.1016/j.tree.2007.06.012)
- 299 7. Sasaki, T., Brännström, Å., Dieckmann, U. & Sigmund, K. 2012 The take-it-or-leave-it
300 option allows small penalties to overcome social dilemmas. *Proc. Natl Acad. Sci. USA*
301 **109**, 1165–1169. (doi:10.1073/pnas.1115219109)
- 302 8. Williams, K. D. 2001 *Ostracism: the power of silence*. New York, NY: Guilford Press.
- 303 9. Masclet, D. 2003 Ostracism in work teams: a public good experiment. *Int. J. Manpow.*
304 **24**, 867–887. (doi:10.1108/01437720310502177)
- 305 10. Cinyabuguma, M., Page, T. & Putterman, L. 2005 Cooperation under the threat of
306 expulsion in a public goods experiment. *J. Public. Econ.* **89**, 1421–1435.
307 (doi:10.1016/j.jpubeco.2004.05.011)
- 308 11. Maier-Rigaud, F. P., Martinsson, P. & Staffiero, G. 2010 Ostracism and the provision
309 of a public good: experimental evidence. *J. Econ. Behav. Organ.* **73**, 387–395.
310 (doi:10.1016/j.jebo.2009.11.001)

- 311 12. Oliver, P. 1980 Rewards and punishments as selective incentives for collective action:
312 Theoretical investigations. *Am. J. Sociol.* **85**, 1356–1375. (doi:10.1086/227168)
- 313 13. Raihani, N. J., Thornton, A. & Bshary, R. 2012 Punishment and cooperation in nature.
314 *Trends. Ecol. Evol.* **27**, 288–295 (doi:10.1016/j.tree.2011.12.004)
- 315 14. Williams, K. D., Cheung, C. K. T. & Choi, W. 2000 Cyberostracism: effects of being
316 ignored over the internet. *J. Pers. Soc. Psychol.* **79**, 748–762. (doi:10.1037/0022-
317 3514.79.5.748)
- 318 15. Kurzban, R. & Leary, M. R. 2001 Evolutionary origins of stigmatization: the functions
319 of social exclusion. *Psychol. Bull.* **127**, 187–208. (doi:10.1037/0033-2909.127.2.187)
- 320 16. Wiessner, P. 2005 Norm Enforcement among the Ju/'hoansi Bushmen. *Hum. Nat.* **16**,
321 115–145. (doi:10.1007/s12110-005-1000-9)
- 322 17. Ostrom, E. 1990 *Governing the commons: the evolution of institutions for collective*
323 *action*. New York, NY: Cambridge University Press.
- 324 18. Cant, M. A. & Johnstone, R. A. 2006 Self-serving punishment and the evolution of
325 cooperation. *Evol. Biol.* **19**, 1383–1385. (doi:10.1111/j.1420-9101.2006.01151.x)
- 326 19. de Waal, F. B. M. & Suchak, M. 2010 Prosocial primates: selfish and unselfish
327 motivations. *Phil. Trans. R. Soc. B* **365**, 2711–2722. (doi:10.1098/rstb.2010.0119)
- 328 20. Melis, A. P., Hare, B. & Tomasello, M. 2006 Chimpanzees recruit the best
329 collaborators. *Science* **311**, 1297–1300. (doi:10.1126/science.1123007)
- 330 21. Bshary R. & Grutter A. S. 2005 Punishment and partner switching cause cooperative
331 behaviour in a cleaning mutualism. *Biol. Lett.* **1**, 396–399.
332 (doi:10.1098/rsbl.2005.0344)
- 333 22. Croson, R., Fatás, E. & Neugebauer, T. 2006 Excludability and contribution: A
334 laboratory study in team production. Working Paper. Wharton School.
- 335 23. Fatas, E., Morales, A. J. & Ubeda, P. 2010 Blind Justice: an experimental analysis of
336 random punishment in team production. *J. Econ. Psychol.* **31**, 358–373.
337 (doi:10.1016/j.joep.2010.01.005)
- 338 24. Ouwerkerk, J. W., Kerr, N. L., Gallucci, M. & Van Lange. P. A. M. 2005 Avoiding the
339 social death penalty: ostracism and cooperation in social dilemmas. In *The social*

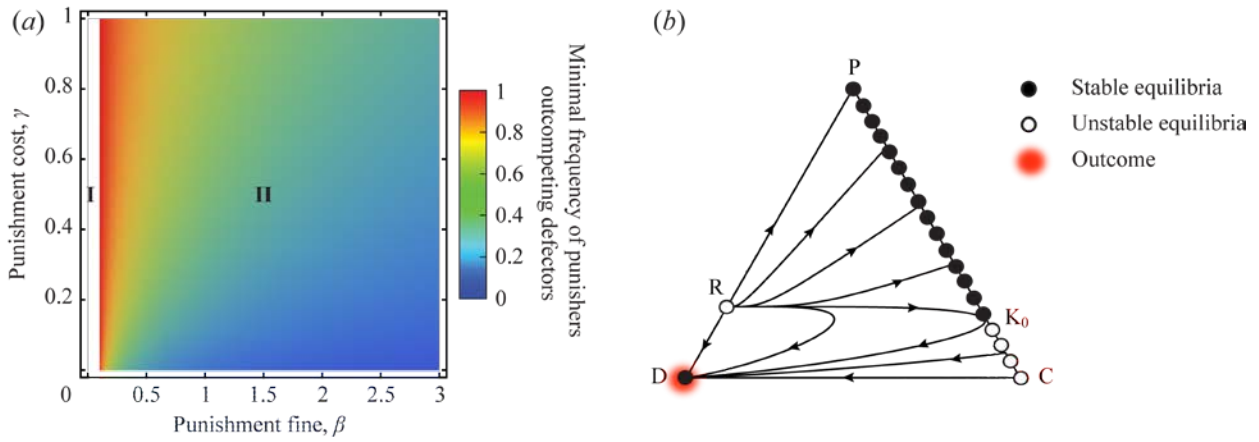
- 340 *outcast: ostracism, social Exclusion, rejection, and bullying* (eds K. D. Williams, J. P.
341 Forgas & W. von Hippel), pp. 321–332. New York, NY: Psychology Press.
- 342 25. Fowler, J. H. 2005 Altruistic punishment and the origin of cooperation. *Proc. Natl*
343 *Acad. Sci. USA* **102**, 7047–7049. (doi:10.1073/pnas.0500938102)
- 344 26. Hauert, C., Traulsen, A., Brandt, H., Nowak, M. A. & Sigmund, K. 2007 Via freedom
345 to coercion: the emergence of costly punishment. *Science* **316**, 1905–1907.
346 (doi:10.1126/science.1141588)
- 347 27. Boyd, R., Gintis, H. & Bowles, S. 2010 Coordinated punishment of defectors sustains
348 cooperation and can proliferate when rare. *Science* **328**, 617–620.
349 (doi:10.1126/science.1183665)
- 350 28. Axelrod, R. 1986 An evolutionary approach to norms. *Am. Polit. Sci. Rev.* **80**, 1095–
351 1111. (doi:10.2307/1960858)
- 352 29. Colman, A. M. 2006 The puzzle of cooperation. *Nature* **440**, 744.
353 (doi:10.1038/440744b)
- 354 30. Hirshleifer, D. & Rasmusen, D. 1989 Cooperation in a repeated prisoners' dilemma
355 with ostracism. *J. Econ. Behav. Organ.* **12**, 87–106. (doi:10.1016/0167-2681(89)90078-
356 4)
- 357 31. Bowls, S. & Gintis, H. 2004 The evolution of strong reciprocity: cooperation in
358 heterogeneous populations. *Theor. Popul. Biol.* **65**, 17–28.
359 (doi:10.1016/j.tpb.2003.07.001)
- 360 32. Panchanathan, K. & Boyd, R. 2004 Indirect reciprocity can stabilize cooperation
361 without the second-order free rider problem. *Nature* **432**, 499–502.
362 (doi:10.1038/nature02978)
- 363 33. Hauert, C., Haiden, N. & Sigmund, K. 2004 The dynamics of public goods. *Discrete*
364 *Continuous Dyn. Syst. Ser. B* **4**, 575–585. (doi:10.3934/dcdsb.2004.4.575)
- 365 34. Hauert, C., Traulsen, A., De Silva née Brandt, H., Nowak, M. A. & Sigmund, K. 2008
366 Public goods with punishment and abstaining in finite and infinite populations. *Biol.*
367 *Theory* **3**, 114–122. (doi:10.1162/biot.2008.3.2.114)

- 368 35. Hofbauer, J. & Sigmund, K. 1998 *Evolutionary Games and Population Dynamics*.
369 Cambridge, UK: Cambridge University Press.
- 370 36. Sigmund, K., De Silva, H., Traulsen, A. & Hauert, C. 2010 Social learning promotes
371 institutions for governing the commons. *Nature* **466**, 861–863.
372 (doi:10.1038/nature09203)
- 373 37. Hilbe, C. & Traulsen, A. 2012 Emergence of responsible sanctions without second
374 order free riders, antisocial punishment or spite. *Scient. Rep.* **2**, 458.
375 (doi:10.1038/srep00458)
- 376 38. Nikiforakis, N. & Normann, H. -T. 2008 A comparative static analysis of punishment in
377 public-good experiment. *Exp. Econ.* **11**, 358–369. (doi:10.1007/s10683-007-9171-3)
- 378 39. Egas, M. & Riedl, A. 2008 The economics of altruistic punishment and the
379 maintenance of cooperation. *Proc. R. Soc. B* **275**, 871–878.
380 (doi:10.1098/rspb.2007.1558)
- 381 40. Sutter, M., Haigner, S. & Kocher, M. G. 2010 Choosing the carrot or the stick?
382 Endogenous institutional choice in social dilemma situations. *Rev. Econ. Stud.* **77**,
383 1540–1566. (doi:10.1111/j.1467-937X.2010.00608.x)
- 384 41. Chuang, J. S., Rivoire, O. & Leibler, S. 2009 Simpson’s paradox in a synthetic
385 microbial system. *Science* **323**, 272–275. (doi:10.1126/science.1166739)
- 386 42. Kiyonari, T. & Barclay, P. 2008 Cooperation in social dilemma: free riding may be
387 thwarted by second-order reward rather than by punishment. *J. Pers. Soc. Psychol.* **95**,
388 826–842. (doi:10.1037/a0011381)
- 389 43. Cornforth, D. M., Sumpter, D. J. T., Brown, S. P. & Brännström, A. 2012 Synergy and
390 group size in microbial cooperation. *Am. Nat.* **180**, 296–305. (doi:10.1086/667193)
- 391 44. Cressman, R., Song, J-W., Zhang, B-Y. & Tao, Y. 2012 Cooperation and evolutionary
392 dynamics in the public goods game with institutional incentives. *J. Theor. Biol.* **299**,
393 144–151. (doi:10.1016/j.jtbi.2011.07.030)
- 394 45. Eshel, I., Samuelson, L. & Shaked, A. 1998 Altruists, egoists and hooligans in a local
395 interaction model. *J. Econ. Theory* **88**, 157–179.

- 396 46. Nowak, M. A., Tarnita, C. E. & Antal, T. 2010 Evolutionary dynamics in structured
397 populations. *Phil. Trans. R. Soc. B* **365**, 19–30. (doi:10.1098/rstb.2009.0215)
- 398 47. Christakis, N. A. & Fowler, J. H. 2012 Social contagion theory: examining dynamic
399 social networks and human behavior. *Stat. Med.* Epub. 2012 Jun 18.
400 (doi:10.1002/sim.5408)
- 401 48. Brandt, H., Hauert, C. & Sigmund, K. 2003 Punishment and reputation in spatial public
402 goods games. *Proc. R. Soc. B* **270**, 1099–1104. (doi:10.1098/rspb.2003.2336)
- 403 49. Nakamaru, M. & Dieckmann, U. 2009 Runaway selection for cooperation and strict-
404 and-severe punishment. *J. Theor. Biol.* **257**, 1–8. (doi:10.1016/j.jtbi.2008.09.004)
- 405 50. Rand, D. G., Armao, J. J., Nakamaru, M. & Ohtsuki, H. 2010 Anti-social punishment
406 can prevent the co-evolution of punishment and cooperation. *J. Theor. Biol.* **265**, 624–
407 632. (doi:10.1016/j.jtbi.2010.06.010)

408

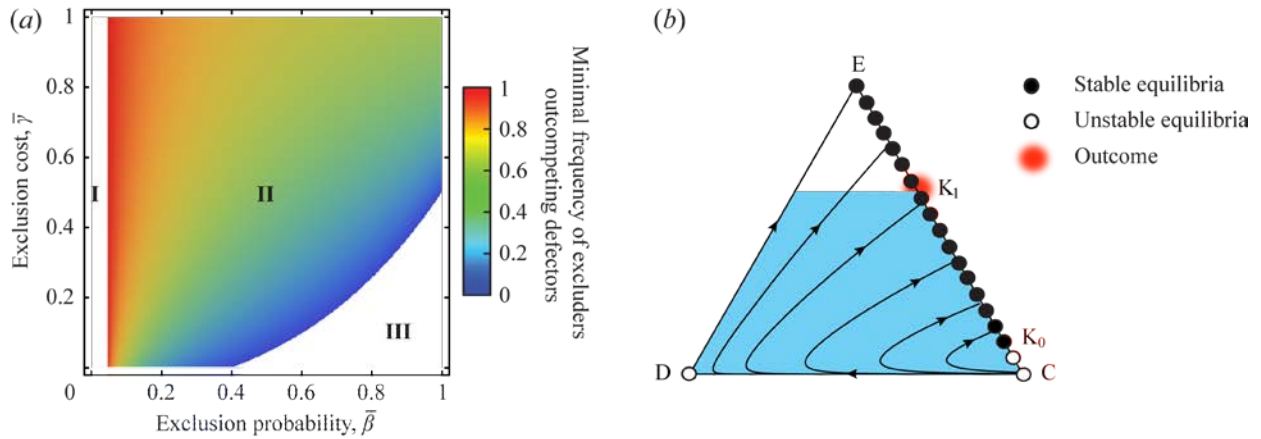
409 **Figure captions**



410

411 Figure 1. Effects of punishing freeriders. (a) Between the punishers and freeriders. **I**, If β is
 412 smaller than a threshold value $\beta_0 = \sigma/(n-1)$, where $\sigma = c(1-r/n)$ describes a net cost for the
 413 single contributor, the defectors dominate. **II**, If β is greater than β_0 , punishing leads to
 414 bistable competition between the two strategies. With increasing β or decreasing γ , the
 415 minimal frequency of the punishers outcompeting the defectors decreases. However, the
 416 excluders cannot dominate the defectors for finitely large values of β . Parameters: group size
 417 $n = 5$, multiplication factor $r = 3$, and contribution cost $c = 1$. (b) In the presence of second-
 418 order freeriders. The triangle represents the state space, $\Delta = \{(x, y, z) : x, y, z \geq 0, x + y + z = 1\}$,
 419 where x, y , and z are the frequencies of the cooperators, defectors, and punishers, respectively.
 420 The vertices, C, D, and P, correspond to the three homogeneous states in which all are the
 421 cooperators ($x = 1$), defectors ($y = 1$), or punishers ($z = 1$). The edge PC consists of a
 422 continuum of equilibria. The defectors dominate the cooperators. Here we specifically assume
 423 $\beta = 0.5$ and $\gamma = 0.03$, which result in an unstable equilibrium R within PD and the
 424 segmentation of PC into stable part PK₀ and unstable part K₀C. The interior of Δ is separated
 425 into the basins of attraction of D and PK₀. In fact, given the occasional mutation to a defector,
 426 the population's state must leave PK₀ and then enter the neighborhood of the unstable segment
 427 K₀C because $P_p > P_c$ holds over the interior space. The population eventually converges to D.

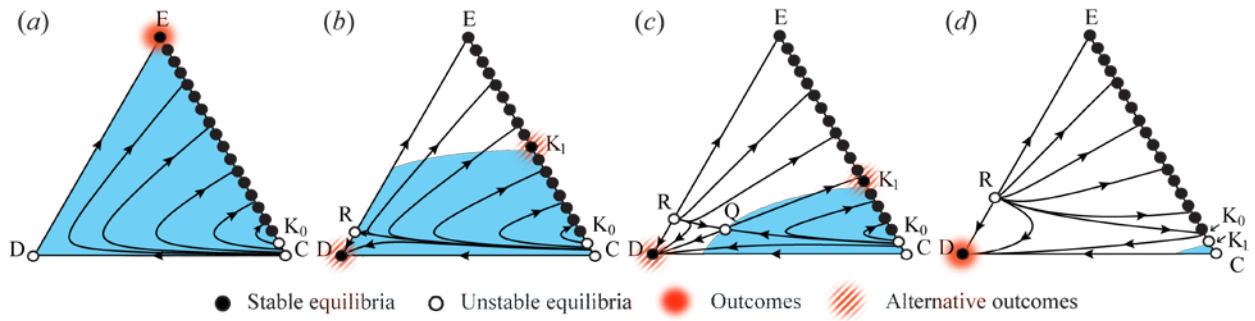
428



429

430 Figure 2. Effects of excluding freeriders. *(a)* Between the excluders and freeriders. **I**, If $\bar{\beta}$ is
 431 smaller than a threshold value z_0 , the defectors dominate. **II**, If $\bar{\beta}$ is greater than z_0 ,
 432 exclusion leads to bistable competition between the two strategies. With increasing $\bar{\beta}$ or
 433 decreasing $\bar{\gamma}$, the minimal frequency of the excluders outcompeting the defectors decreases.
 434 **III**, If $\bar{\beta}$ and $\bar{\gamma}$ are sufficiently high and low, the excluders dominate. The parameters are as
 435 in figure 1*a*. *(b)* In the presence of second-order freeriders. The triangle Δ is as in figure 1*b*,
 436 except that z denotes the excluder frequency and the vertex E corresponds to its homogeneous
 437 state. Similarly, the edge EC consists of a continuum of equilibria. Here we specifically
 438 assume $\bar{\beta} = 1$ and $\bar{\gamma} = 0.03$. EC is separated into stable and unstable segments. The coloured
 439 area in the interior of Δ is the region in which $P_E > P_C$ holds. In fact, given the occasional
 440 mutation to a defector, the population's state must converge to the vicinity of the point K₁,
 441 because the advantage of the excluders over the cooperators becomes broken when the
 442 population's state goes up beyond K₁.

443



444

445 Figure 3. Effects of intermediate social exclusion in the presence of second-order freeriders.
 446 The parameters and triangles are as in figure 1, except that $\bar{\beta} = 0.5$ and $\bar{\gamma} = 0.03$ (a), 0.13 (b),
 447 0.18 (c), or 0.28 (d). EC is separated into stable and unstable segments. The coloured area is
 448 the interior region in which $P_E > P_C$ holds. (a) The dynamics of ED are unidirectional to E. All
 449 interior trajectories converge onto the stable segment EK_0 . Moreover, occasionally mutating to
 450 a defector leads to upgrading E to a global attractor. (b-d) An unstable equilibrium R appears
 451 on ED. The interior space is separated into the basins of attraction of D and EK_0 . R is a saddle
 452 (b) or source (c and d). In (c) especially, the interior space has a saddle point Q. Given the
 453 mutant defectors, the population's state around EK_0 will gradually move to K_1 (b and c), or to
 454 the unstable segment K_0C (d). The last case is followed by a convergence toward D.

Electronic supplementary material (ESM)

This includes: Materials and methods, and Supplementary figures, S1–S6

Materials and methods

We first determine the strategy's payoffs in public good games with social exclusion, then show details of individual-based simulations for assessing the robustness with respect to stochastic evolutionary game dynamics.

Payoffs for social exclusion: We consider the replicator dynamics for the cooperator (C), defector (D), and excluder (E), with frequencies of x , y , and z , respectively. Thus, $x, y, z \geq 0$ and $x + y + z = 1$. We denote the expected payoff values for the three strategies by P_S , with $S = C, D, \text{ and } E$, respectively. The replicator system is given by

$$\dot{x} = x(P_C - \bar{P}), \quad \dot{y} = y(P_D - \bar{P}), \quad \dot{z} = z(P_E - \bar{P}),$$

where $\bar{P} := xP_C + yP_D + zP_E$ describes the average payoff in the entire population. We denote by X, Y , and Z the number of the cooperators, defectors, and excluders, respectively, among the $(n-1)$ coplayers around a focal player. Then, if W of the Y defectors have not been excluded by every excluder, the expected payoff for each strategy is given by

$$P_S = \sum_{X=0}^{n-1} \sum_{Y=0}^{n-1-X} \sum_{W=0}^Y \pi_S p_S. \quad (\text{S1})$$

In equation (S1), p_S denotes the payoff for the focal player who follows the strategy S among the $(n-1)$ coplayers with a configuration of $\{X, Y, Z, W\}$, and π_S denotes the probability to find the specified coplayers. Using a function $\alpha(Z)$ that denotes the probability that all of the Z excluders fail to exclude a targeted defector, we have

$$p_C = \frac{rc(X+Z+1)}{X+W+Z+1} - c, \quad (\text{S2})$$

$$p_D = \alpha(Z) \frac{rc(X+Z)}{X+W+Z+1}, \quad (\text{S3})$$

$$p_E = p_C - \bar{\gamma}Y, \quad (\text{S4})$$

24
$$\pi_C = \pi_D = \binom{n-1}{X, Y, Z} x^X y^Y z^Z \binom{Y}{W} \alpha(Z)^W [1 - \alpha(Z)]^{Y-W}, \quad (\text{S5})$$

25
$$\pi_E = \binom{n-1}{X, Y, Z} x^X y^Y z^Z \binom{Y}{W} \alpha(Z+1)^W [1 - \alpha(Z+1)]^{Y-W}. \quad (\text{S6})$$

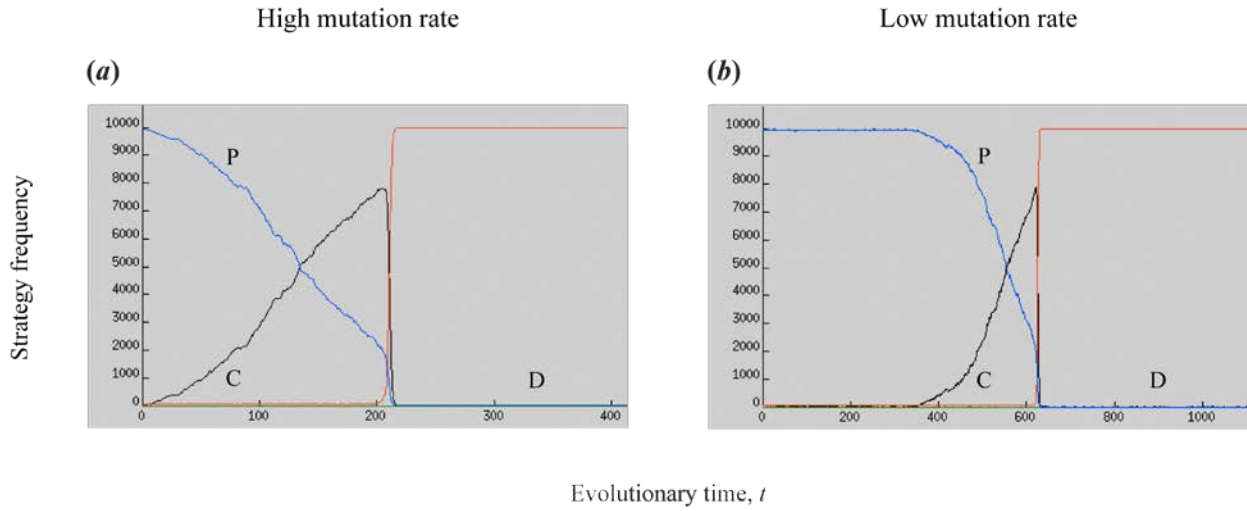
26 In equations (S5) and (S6), $\binom{n-1}{X, Y, Z}$ and $\binom{Y}{W}$ represent the multinomial and binomial
 27 coefficients. Thus, $\binom{n-1}{X, Y, Z} x^X y^Y z^Z$ describes the probability of finding the $(n-1)$ coplayers
 28 with X cooperators, Y defectors, and Z excluders, and $\binom{Y}{W} \alpha(Z)^W [1 - \alpha(Z)]^{Y-W}$ describes the
 29 probability that W of the Y defectors have not been excluded. In the paper, we assume
 30 $\alpha(Z) = (1 - \bar{\beta})^Z$, where $\bar{\beta}$ is the exclusion probability: an excluder succeeds in excluding a
 31 defector.

32 **Individual-based simulation:** Here, we consider a finitely large, well-mixed population with
 33 M interacting individuals. For the dynamic analysis, instead of the replicator system [35], we
 34 implement a pairwise comparison process among finite individuals [36,37], which is based on
 35 preferentially imitating strategies of more successful individuals. We assume that the
 36 individual strategies are updated asynchronously as follows. First, an individual i is selected at
 37 random and then earns its “average” payoff p_i after engaging in T games with coplayers
 38 randomly selected in each case. Second, the focal individual i faces a model individual j who is
 39 drawn at random, with its average payoff p_j that is calculated throughout independent T
 40 games. If $p_i \geq p_j$, no update occurs; or otherwise, i will adopt j 's strategy, with the probability
 41 given by

42
$$\theta_{i \rightarrow j} = \frac{1}{1 + \exp(-K(p_j - p_i))},$$

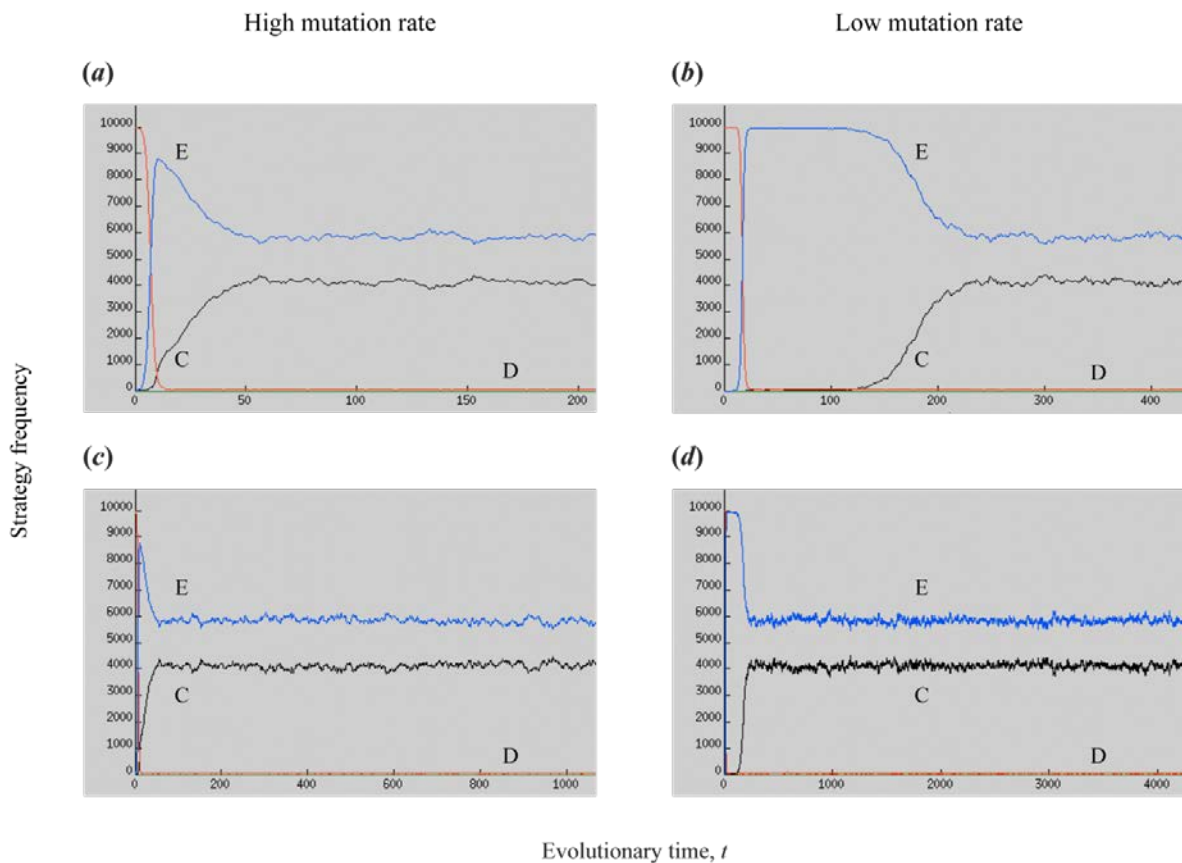
43 where K denotes the selection strength. Finally, the focal individual i can mutate and turn into a
 44 cooperator, defector, or punisher (or excluder) with probabilities μ_C, μ_D, μ_P (or μ_E). Our
 45 numerical results demonstrated in figures S1–S6 are robust with respect to changes in the
 46 parameter values of $M, \mu_C, \mu_D, \mu_P, \mu_E$, and K .

47 **Supplementary figures**



48

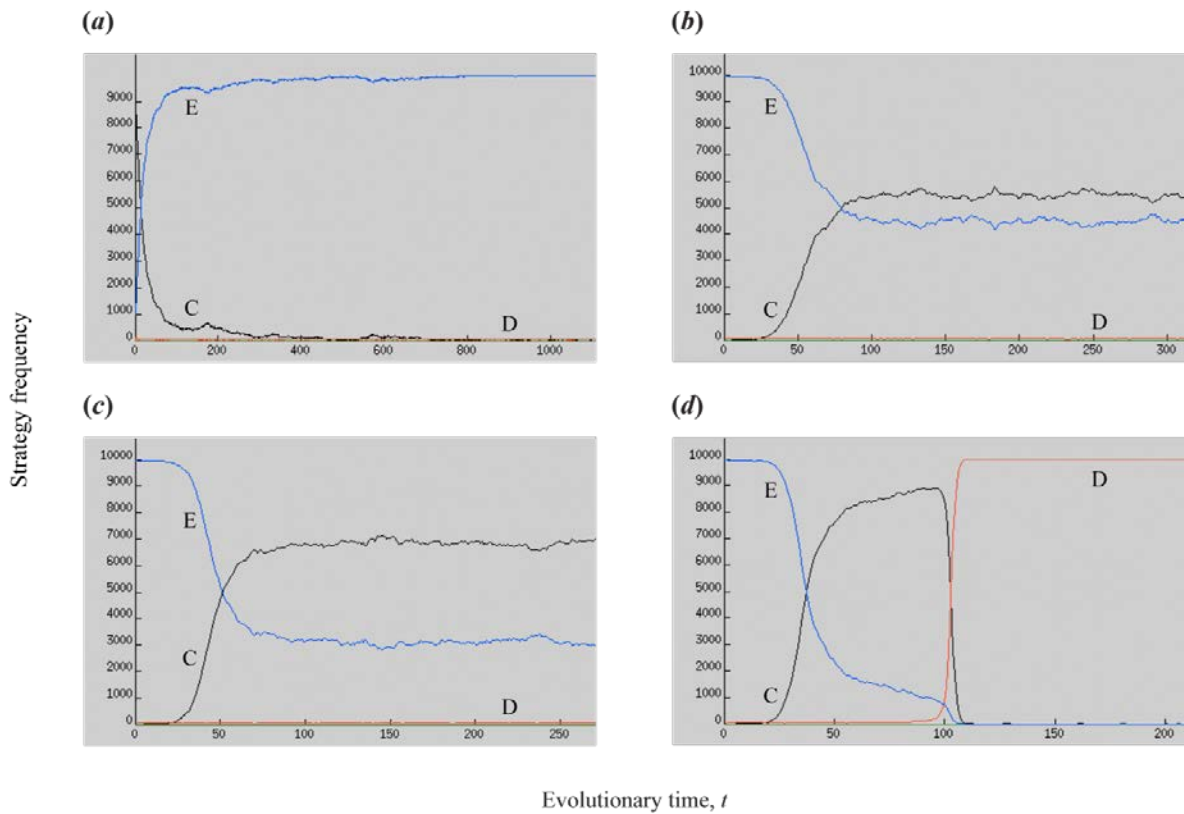
49 Figure S1. Individual-based simulation for public good games with costly punishment. We
 50 began with a 100%-punisher population to observe its stability. First, because the punishing of
 51 mutant defectors is costly, the former major punishers (blue) will gradually be replaced by the
 52 initially minor cooperators (namely, second-order freeriders, black). Next, when a critical
 53 fraction of punishers is lost, the mutant defectors (red) succeed in invading the population and
 54 then quickly prevail. The parameters are as in figure 1b: group size $n = 5$, multiplication factor
 55 $r = 3$, contribution cost $c = 1$, punishment cost $\beta = 0.5$, and punishment fine $\gamma = 0.03$. The
 56 defectors dominate the cooperators, and the excluders and defectors are under bistable
 57 competition. Other parameters are as the population size $M = 10^4$, sample game count $T = 50$,
 58 selection strength $K = 200$, mutation rate to D $\mu_D = 5 \times 10^{-3}$, mutation rates to C and P
 59 $\mu_C = \mu_P = 10^{-5}$ (low mutation rate) or $\mu_C = \mu_P = 10^{-3}$ (high mutation rate), and the unit of
 60 evolutionary time t describes 10^4 times the iteration of the update sequence.



61

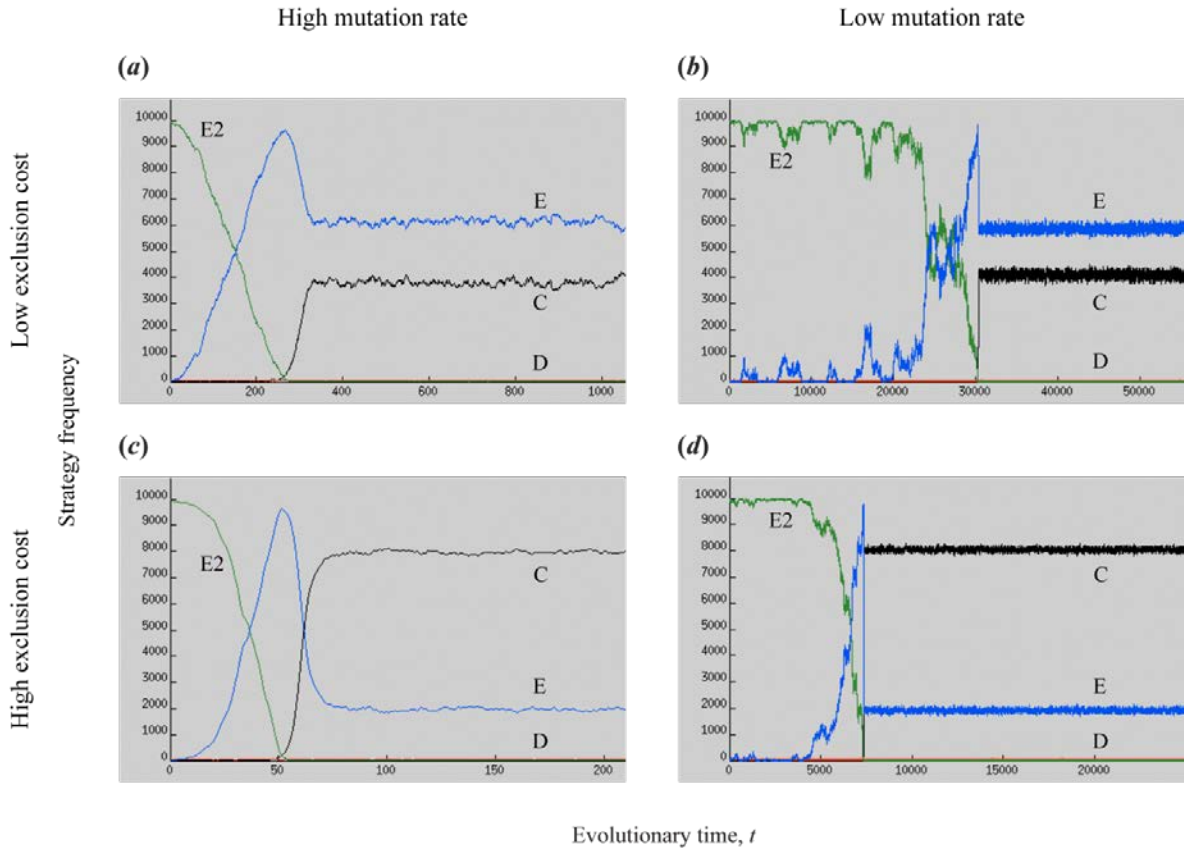
Evolutionary time, t

62 Figure S2. Individual-based simulation for public good games with perfect social exclusion.
 63 The parameters are as in figure 2b: $n = 5$, $r = 3$, $c = 1$, exclusion probability $\bar{\beta} = 1$, and
 64 exclusion cost $\bar{\gamma} = 0.03$. We began with a 100%-punisher population to observe the
 65 establishment of a cooperative state. Whether the minimal mutation rate is high (10^{-3}) or low
 66 (10^{-5}), the former major defectors (red) will soon be replaced by the initially minor excluders
 67 (blue), whose part will then be gradually replaced by the cooperators (black). The population
 68 eventually converges to a certain mixture state of the contributors without a second-order
 69 freerider problem. The final state has been indicated by point K_1 in figures 2b. The simulation
 70 parameters are as in figure S1.



71 Evolutionary time, t

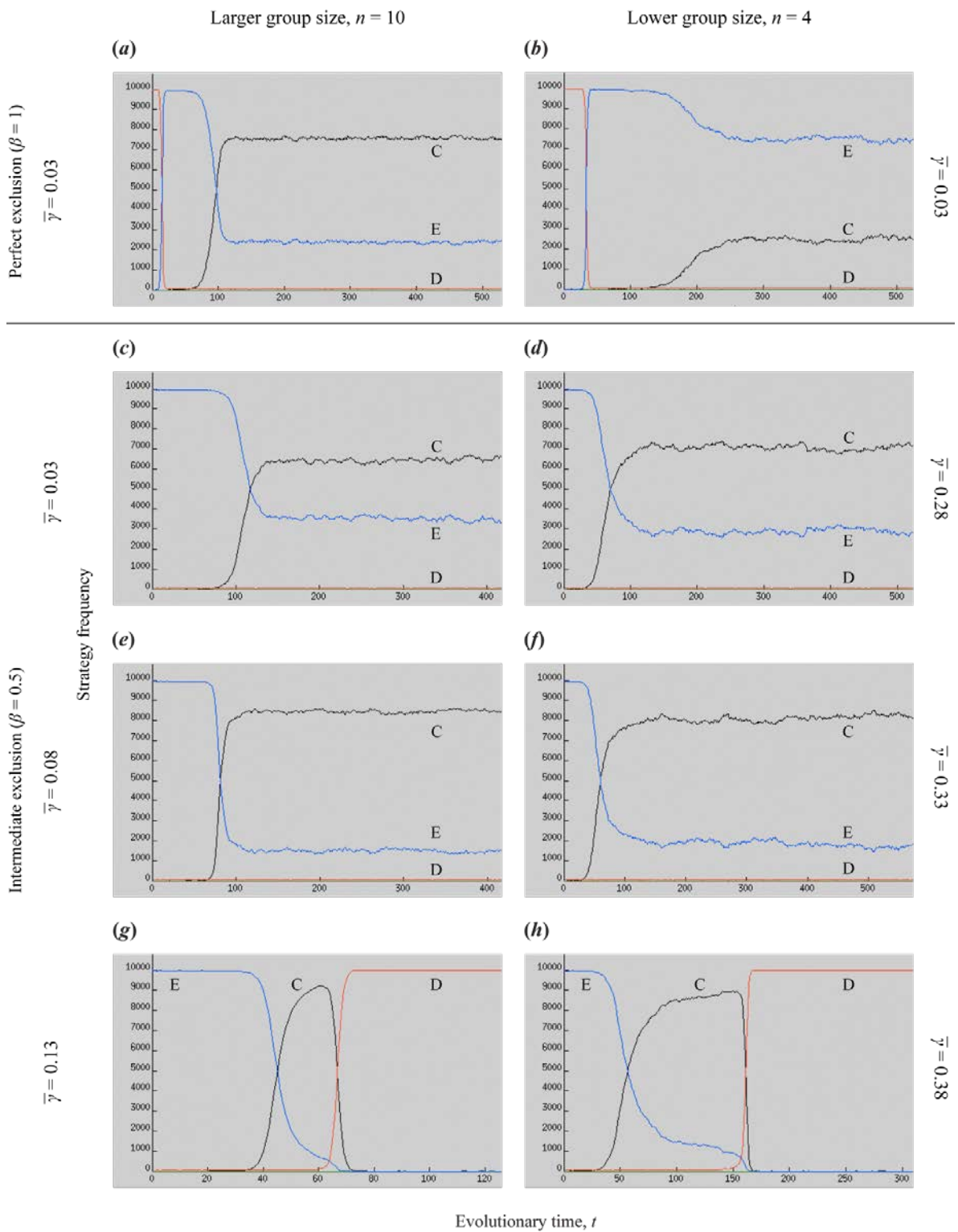
72 Figure S3. Individual-based simulation for public good games with intermediate social
 73 exclusion. The parameters are as in figure 3: $n = 5$, $r = 3$, $c = 1$, and $\bar{\beta} = 0.5$. We began with
 74 different initial conditions, depending on the value of $\bar{\gamma}$: 90% cooperators and 10% excluders
 75 for $\bar{\gamma} = 0.03$ (a) and 100% excluders for $\bar{\gamma} = 0.13$ (b), 0.18 (c), or 0.28 (d). (a) The former
 76 major cooperators (black) will gradually be replaced by the initially minor excluders (blue),
 77 which then stably occupy the entire population (b and c). The initially minor cooperators will
 78 first replace part of the excluders, and the population will then converge to a certain mixture
 79 state, which has been indicated by the point K_1 in figures 3b and 3c, respectively (d). As in (b
 80 and c), the cooperators will gradually expand. When a critical fraction of the excluders is lost
 81 (the point K_0), the mutant defectors (black) succeed in invading the population and will then
 82 quickly prevail to 100%. The simulation parameters are as in figure S1 with the low mutation
 83 rate.



84

Evolutionary time, t

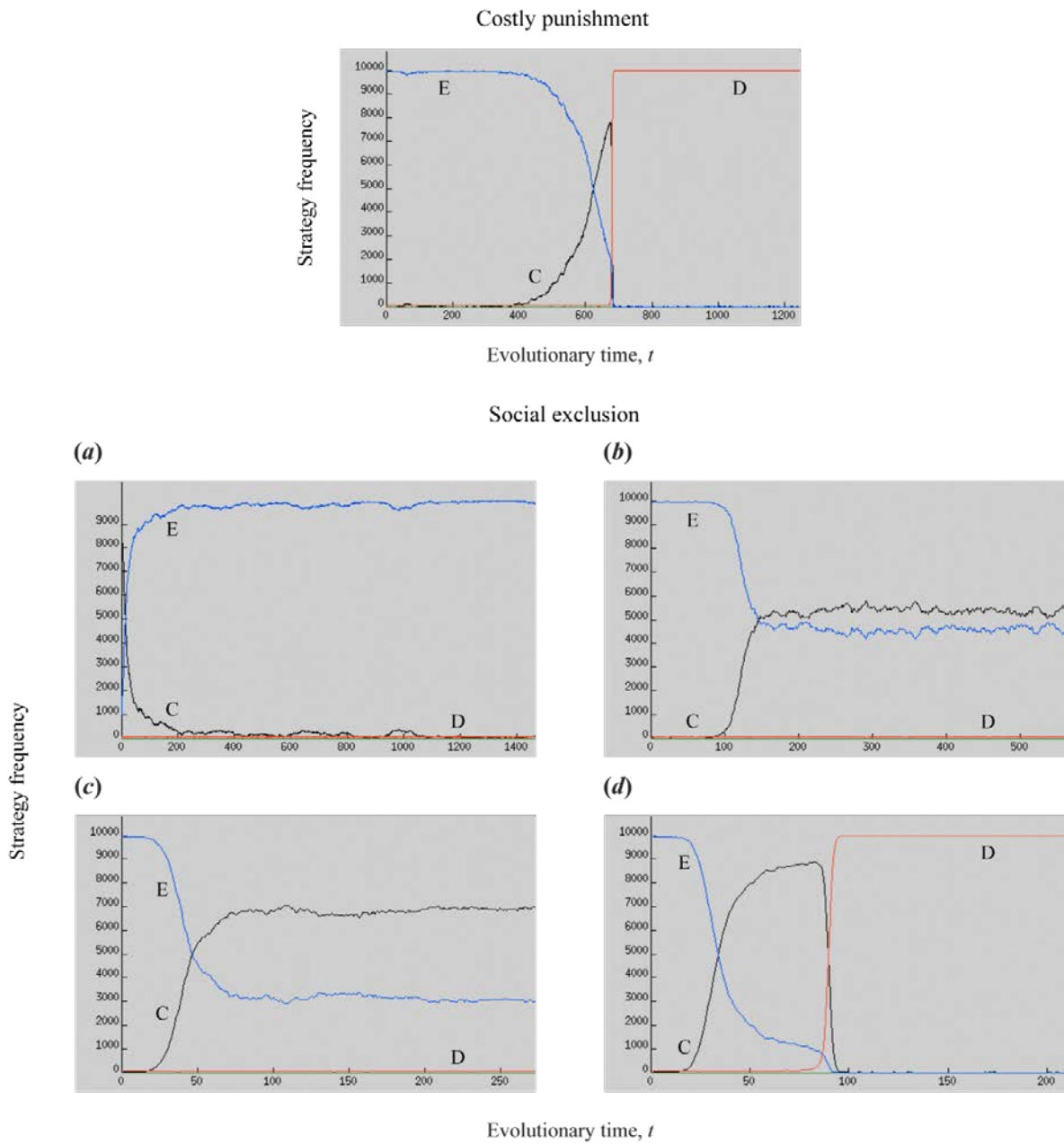
85 Figure S4. Individual-based simulation for public good games with second-order social
 86 exclusion. The parameters are as in figure 2b, except that $\bar{\gamma} = 0.03$ (low exclusion cost) or
 87 $\bar{\gamma} = 0.28$ (high exclusion cost). We began with the initial condition: 100% second-order
 88 excluders (green) who in the presence of the defectors, also exclude the cooperators, as well as
 89 the defectors (with the same cost and probability). The initial residents will first be replaced
 90 with the excluders (blue), and then are partially invaded by the cooperators (black): the
 91 population will converge to a certain mixture state of the contributors, whether with a high or
 92 low exclusion cost. The simulation parameters are as in figure S1.



93

Evolutionary time, t

94 Figure S5. Effect of different group sizes. The parameters are as in figure 2b, for perfect
 95 exclusion (a) and (b), and in figure 3, for intermediate exclusion (c–h). The initial conditions
 96 are 100% second-order excluders in (a) and (b) and 100% excluders in (c–h).



97

98 Figure S6. Effect of options to choose the number of sanctioned defectors. The model and
 99 simulation parameters, and initial conditions are as in figure S1, for costly punishment (top),
 100 and in figure S3, for intermediate exclusion (middle and bottom, *a-d*). Here we assume that a
 101 punisher or excluder is willing to sanction only one defector selected at random from all
 102 defectors in the group. The results are almost same as in figures S1 and S3.